

# EFFICIENT METHOD FOR BI-LEVEL OPTIMIZATION WITH NON-SMOOTH LOWER-LEVEL PROBLEM

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Bi-level optimization plays a key role in a lot of machine learning applications. Existing state-of-the-art bi-level optimization methods are limited to smooth or some specific non-smooth lower-level problems. Therefore, achieving an efficient algorithm for the bi-level problems with a generalized non-smooth lower-level objective is still an open problem. To address this problem, in this paper, we propose a new bi-level optimization algorithm based on smoothing and penalty techniques. Using the theory of generalized directional derivative, we derive new conditions for the bilevel optimization problem with nonsmooth, perhaps non-Lipschitz lower-level problem, and prove our method can converge to the points satisfying these conditions. We also compare our method with existing state-of-the-art bi-level optimization methods and demonstrate that our method is superior to the others in terms of accuracy and efficiency.

## 1 INTRODUCTION

Bi-level optimization (BO) (Bard, 2013; Colson et al., 2007) plays a central role in various machine learning applications including hyper-parameter optimization (Pedregosa, 2016; Bergstra et al., 2011; Bertsekas, 1976), meta-learning (Feurer et al., 2015; Franceschi et al., 2018; Rajeswaran et al., 2019), reinforcement learning (Hong et al., 2020; Konda & Tsitsiklis, 2000). It involves a competition between two parties or two objectives, and if one party makes its choice first it will affect the optimal choice for the other party. Several approaches, such as Bayesian optimization (Klein et al., 2017), random search (Bergstra & Bengio, 2012), evolution strategy (Sinha et al., 2017), gradient-based methods (Pedregosa, 2016; Maclaurin et al., 2015; Swersky et al., 2014), have been proposed to solve BO problems, among which gradient-based methods have become the mainstream for large-scale BO problems.

The key idea of the gradient-based method is to approximate the gradient of upper-level variables, called hypergradient. For example, the implicit differentiation methods (Pedregosa, 2016; Rajeswaran et al., 2019) use the first derivative of the lower-level problem to be 0 to derive the hypergradient. The explicit differentiation methods calculate the gradient of the update rules of the lower-level based on chain rule (Maclaurin et al., 2015; Domke, 2012; Franceschi et al., 2017; Swersky et al., 2014) to approximate the hypergradient. Mehra & Hamm (2019) reformulate the bilevel problem as a single-level constrained problem by replacing the lower level problem with its first-order necessary conditions, and then solve the new problem by using the penalty method. Obviously, all these methods need the lower-level problem to be smooth.

However, in many real-world applications, such as image restoration (Chen et al.; Nikolova et al., 2008), variable selection (Fan & Li, 2001; Huang et al., 2008; Zhang et al., 2010) and signal processing (Bruckstein et al., 2009), the objective may have a complicated non-smooth, perhaps non-Lipschitz term (Bian & Chen, 2017). Traditional methods cannot be directly used to solve the bilevel problem with such a lower-level problem. To solve the BO problems with some specific nonsmooth lower-level problems, researchers have proposed several algorithms based on the above-mentioned methods. Specifically, Bertrand et al. (2020) searched the regularization parameters for LASSO-type problems by approximating the hypergradient from the soft thresholding function (Donoho, 1995; Bredies & Lorenz, 2008; Beck & Teboulle, 2009). Frecon et al. (2018) proposed a primal-dual FMD-based method, called FBBGLasso, to search the group structures of group-LASSO problems. Okuno et al. (2021) used the smoothing method and constrained optimization method to search the regularization

Table 1: Representative gradient-based bi-level optimization methods.

Method	Reference	Problem	Method type
FMD	Franceschi et al. (2017)	Smooth	Bi-level
RMD	Franceschi et al. (2017)	Smooth	Bi-level
Approx	Pedregosa (2016)	Smooth	Bi-level
Penalty	Mehra & Hamm (2019)	Smooth	Single-level
FBGGL	Frecon et al. (2018)	Group LASSO	Bi-level
SparseHO	Bertrand et al. (2020)	LASSO-type	Bi-level
SMNBP	Okuno et al. (2021)	$p$ -norm	Single-level
SPNBO	Ours	Generalized	Single-level

parameter of  $q$ -norm ( $0 < q \leq 1$ ) and provided the convergence analysis of their method. We summarize several representative methods in Table 1. Obviously, all these methods and their theoretic analysis only focus on some specific problem and can not be used to solve the bilevel problem with a generalized nonsmoothed lower-level problem. Therefore, how to solve the BO problem with a generalized non-smooth lower-level objective and obtain its convergence analysis are still open problems.

To address this problem, in this paper, we propose a new algorithm, called SPNBO, based on smoothing (Nesterov, 2005; Chen et al., 2013) and penalty (Wright & Nocedal, 1999) techniques. Specifically, we use the smoothing technique to approximate the original non-Lipschitz lower-level problem and generate a sequence of smoothed bi-level problems. Then, a single-level constrained problem is obtained by replacing the smoothed lower-level objective with its first-order necessary condition. For each given smoothing parameter, we propose a stochastic constraint optimization method to solve the single-level constrained problem to avoid calculating the Hessian matrix of the lower-level problem. Theoretically, using the theory of generalized directional derivative, we derive new conditions for the bilevel optimization problem with nonsmooth, perhaps non-Lipschitz lower-level problem, and prove our method can converge to the points satisfying these conditions. We also compare our method with several state-of-the-art bi-level optimization methods, and the experimental results demonstrate that our method is superior to the others in terms of accuracy and efficiency.

**Contributions.** We summarize the main contributions of this paper as follows:

1. We propose a new method to solve the non-Lipschitz bilevel optimization problem based on the penalty method and smoothing method. By using the stochastic constraint method, our method can avoid calculating the Hessian matrix of the lower-level problem, which makes our method a lower time complexity.
2. Based on the Clarke generalized directional derivative, we propose new conditions for the bilevel problem with a generalized non-smoothed lower-level problem. We prove that our method can converge to the proposed conditions.

## 2 PRELIMINARIES

### 2.1 FORMULATION OF NON-SMOOTH BI-LEVEL OPTIMIZATION PROBLEM

In this paper, we consider the following non-smooth bi-level optimization problem:

$$\min_{\lambda} f(\mathbf{w}^*, \lambda) \quad s.t. \quad \mathbf{w}^* \in \arg \min_{\mathbf{w}} g(\mathbf{w}, \bar{\lambda}) + \exp(\lambda_1) \varphi(h(\mathbf{w})), \quad (1)$$

where  $\lambda := [\lambda_1, \lambda_2, \dots, \lambda_m]^T \in \mathbb{R}^m$ ,  $\bar{\lambda} := [\lambda_2, \dots, \lambda_m]^T$  and  $\mathbf{w} \in \mathbb{R}^d$ .  $f : \mathbb{R}^d \times \mathbb{R}^m \mapsto \mathbb{R}$  and  $g : \mathbb{R}^d \times \mathbb{R}^m \mapsto \mathbb{R}$  are twice continuously differentiable on  $\mathbf{w}$  and  $\lambda$ .  $\varphi(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}$  is twice continuously differentiable.  $h(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}^n$  is continuous, not necessarily convex, not differentiable, or even not Lipschitz at some points. Assume  $h(\mathbf{w}) := (h_1(\mathbf{D}_1^T \mathbf{w}), h_2(\mathbf{D}_2^T \mathbf{w}), \dots, h_n(\mathbf{D}_n^T \mathbf{w}))$ , where  $\mathbf{D}_i \in \mathbb{R}^{d \times r}$  and  $h_i : \mathbb{R}^d \mapsto \mathbb{R}$  ( $i = 1, 2, \dots, n$ ) is continuous. For a fixed point  $\bar{\mathbf{w}}$ , assume we have an index set  $\mathcal{I}_{\bar{\mathbf{w}}} = \{i \in \{1, 2, \dots, n\} : h_i \text{ is not Lipschitz continuous at } \mathbf{D}_i^T \bar{\mathbf{w}}\}$  and if  $i \notin \mathcal{I}_{\bar{\mathbf{w}}}$ ,  $h_i$  is twice continuously differentiable.

## 2.2 EXAMPLES OF NON-SMOOTH NON-LIPSCHITZ LOWER-LEVEL PROBLEMS

The non-smooth non-Lipschitz optimization problems widely exist in image restoration (Chen et al.; Nikolova et al., 2008), variable selection (Fan & Li, 2001; Huang et al., 2008; Zhang et al., 2010) and signal processing (Bruckstein et al., 2009). Here, we give two examples as follows.

1.  $l_p$ -norm (Chen et al., 2013):  $\min_{\mathbf{w}} g(\mathbf{w}, \bar{\lambda}) + \exp(\lambda_1) \sum_{i=1}^d |\mathbf{w}_i|^p$ , where  $p \in (0, 1]$ .
2. OSCAR penalty (Bondell & Reich, 2008):  $\min_{\mathbf{w}} g(\mathbf{w}, \bar{\lambda}) + \exp(\hat{\lambda}) \|\mathbf{w}\|_1 + \exp(\tilde{\lambda}) \sum_{i < j} \max\{\mathbf{w}_{\mathcal{G}_i}, \mathbf{w}_{\mathcal{G}_j}\}$ , where  $\mathcal{G}_i$  denotes the group index.

Note that Okuno et al. (2021) only considered the bilevel problem with the lower-level problem given in Example 1. Their theoretical analysis is not suitable for the problem in Example 2 or even more complicated formulation.

## 3 PROPOSED METHOD

In this section, we give a brief review of the smoothing method and then propose our stochastic gradient algorithm based on the penalty method and single-level reduction method to solve the bilevel problem.

### 3.1 SMOOTHING TECHNIQUE

Here, we give the definition of smoothing function (Nesterov, 2005; Chen et al., 2013; Bian & Chen, 2017) which is widely used in nonsmooth non-Lipschitz problems.

**Definition 1.** Let  $\psi : \mathbb{R}^d \mapsto \mathbb{R}$  be a continuous nonsmooth, non-Lipschitz function. We call  $\tilde{\psi} : \mathbb{R}^d \times [0, +\infty) \mapsto \mathbb{R}$  a smoothing function of  $\psi$ , if  $\tilde{\psi}(\cdot, \mu)$  is twice continuously differentiable for any fixed  $\mu > 0$  and  $\lim_{\hat{\mathbf{w}} \mapsto \mathbf{w}, \mu \rightarrow 0} \tilde{\psi}(\hat{\mathbf{w}}, \mu) = \psi(\mathbf{w})$  holds for any  $\mathbf{w} \in \mathbb{R}^d$ .

Here, we give two examples of smoothing functions. The smoothing function of  $\psi_1(w) = \sum_{i=1}^d |w_i|$  is  $\tilde{\psi}_1(w, \mu) = \sum_{i=1}^d (w_i^2 + \mu^2)^{1/2}$  and the smoothing function of  $\psi_2(w) = \sum_{i < j} \max\{w_i, w_j\}$  is  $\tilde{\psi}_2(w, \mu) = \sum_{i < j} \frac{1}{2} (\sqrt{(w_i + w_j)^2 + \mu^2} + \sqrt{(w_i - w_j)^2 + \mu^2})$ .

According to Definition 1, the non-smooth lower level problem in problem (1) could be approximated by using a sequence of the following parameterized smoothing functions,

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} g(\mathbf{w}, \bar{\lambda}) + \exp(\lambda_1) \varphi(\tilde{h}(\mathbf{w}, \mu^k)) \quad (2)$$

where  $\mu^k > 0$  is the smoothing parameter and  $\tilde{h}(\mathbf{w}, \mu^k) := (\tilde{h}_1(\mathbf{D}_1^T \mathbf{w}, \mu^k), \tilde{h}_2(\mathbf{D}_2^T \mathbf{w}, \mu^k), \dots, \tilde{h}_n(\mathbf{D}_n^T \mathbf{w}, \mu^k))$ .

For each given smoothing parameter  $\mu^k > 0$ , we can replace the smoothed lower-level objective with its first-order necessary condition and derive the following single-level problem:

$$\min_{\mathbf{w}, \lambda} f(\mathbf{w}, \lambda) \quad s.t. \quad c(\mathbf{w}, \lambda; \mu^k) = \mathbf{0}, \quad (3)$$

where  $c(\mathbf{w}, \lambda; \mu^k) := \nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\lambda}) + \exp(\lambda_1) \nabla_{\mathbf{w}} \varphi(\tilde{h}(\mathbf{w}, \mu^k))$  and  $\nabla_{\mathbf{w}} \varphi(\tilde{h}(\mathbf{w}, \mu^k)) = \varphi'(z)_{z=\tilde{h}(\mathbf{w}, \mu^k)} \nabla_{\mathbf{w}} \tilde{h}(\mathbf{w}, \mu^k)$ .

### 3.2 STOCHASTIC CONSTRAINT GRADIENT METHOD

In this subsection, we discuss our method to solve the subproblem (3). Obviously, we can use the gradient method to solve its corresponding penalty function to solve the single-level constrained problem. However, calculating the gradient of the penalty functions needs to calculate the Hessian matrix. If the dimension of  $\mathbf{w}$ , calculating the Hessian matrix is very time-consuming. To solve this problem, we introduce a stochastic layer into the constraint such that we only need to calculate the

**Algorithm 1** Smoothing and Penalty Method for Non-Lipschitz Bi-level Optimization (SPNBO)**Input:**  $K, \mu^1, \beta^1, \delta_\mu, \delta_\epsilon \in (0, 1)$ .**Output:**  $\mathbf{w}^{k+1}$  and  $\boldsymbol{\lambda}^{k+1}$ .

```

1: for  $k = 1, \dots, K$  do
2:   Find  $(\mathbf{w}^{k+1}, \boldsymbol{\lambda}^{k+1}, \mathbf{p}^{k+1}) := \min_{\mathbf{w}, \boldsymbol{\lambda}} \max_{\mathbf{p} \in \Delta^d} \mathcal{L}(\mathbf{w}^k, \boldsymbol{\lambda}^k, \mathbf{p}^k, \mu^k)$  using the SCG method.
3:    $\mu^{k+1} = \delta_\mu \mu^k$ .
4:    $\epsilon^{k+1} = \delta_\epsilon \epsilon^k$ .
5: end for

```

gradient of the sampled element of the constraint. Specifically, we reformulate the subproblem (3) as the following minimax problem

$$\min_{\mathbf{w}, \boldsymbol{\lambda}} \max_{\mathbf{p} \in \Delta^d} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k) = f(\mathbf{w}, \boldsymbol{\lambda}) + \beta \sum_{i=1}^d p_i c_i^2(\mathbf{w}, \boldsymbol{\lambda}; \mu^k) - \frac{\tau}{2} \|\mathbf{p}\|_2^2, \quad (4)$$

where  $\beta > 0$ ,  $\lambda > 0$ ,  $\mathbf{p} \in \Delta^d := \{\mathbf{p} | \sum_{i=1}^d p_i = 1 \& 0 \leq p_i \leq 1\}$ ,  $c_i(\mathbf{w}, \boldsymbol{\lambda}; \mu^k)$  denote the  $i$ -th elements of  $c(\mathbf{w}, \boldsymbol{\lambda}; \mu^k)$ . The last term is used to ensure  $\mathcal{L}$  is strongly-concave on  $\mathbf{p}$ . Such a reformulation is widely used in many methods (Cotter et al., 2016; Narasimhan et al., 2020; Shi et al., 2022) to solve the constrained problem.

In each iteration, we sample an element  $w_i$  of  $\mathbf{w}$  according to distribution  $p$  and calculate the corresponding value of  $c_i$  and its gradient w.r.t  $\mathbf{w}$ . Then, we can obtain the stochastic gradient of  $\mathcal{L}$  w.r.t  $\mathbf{w}$  as follows,

$$\hat{\nabla}_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t) = \nabla_{\mathbf{w}} f(\mathbf{w}_t, \boldsymbol{\lambda}_t) + 2\beta c_i(\mathbf{w}_t, \boldsymbol{\lambda}_t; \mu^k) \nabla_{\mathbf{w}} c_i(\mathbf{w}_t, \boldsymbol{\lambda}_t; \mu^k). \quad (5)$$

Using the same method, we can obtain the stochastic gradient  $\hat{\nabla}_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t)$ . Then,

**Algorithm 2** Stochastic constraint gradient (SCG)**Input:**  $\gamma, \sigma, \eta_t, a_{t+1,1}, a_{t+1,2}$ **Output:**  $\mathbf{w}$  and  $\boldsymbol{\lambda}$ .

```

1: Initialize  $m_{1,1}, m_{1,2}, m_{1,3}, \hat{m}_{t,1}, \hat{m}_{t,2}, \hat{m}_{t,3}, \eta_1$ .
2: while Not satisfy the conditions (10) do
3:    $\tilde{\mathbf{w}}_{t+1} = \mathbf{w}_t - \gamma A_{t,1}^{-1} m_{t,1}$ .
4:    $\mathbf{w}_{t+1} = \mathbf{w}_t + \eta_t (\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t)$ .
5:    $\tilde{\boldsymbol{\lambda}}_{t+1} = \boldsymbol{\lambda}_t - \gamma A_{t,2}^{-1} m_{t,2}$ .
6:    $\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t + \eta_t (\tilde{\boldsymbol{\lambda}}_{t+1} - \boldsymbol{\lambda}_t)$ .
7:    $\tilde{\mathbf{p}}_{t+1} = \mathcal{P}_\Delta (\mathbf{p}_t + \sigma A_{t,3}^{-1} m_{t,3})$ .
8:    $\mathbf{p}_{t+1} = \mathbf{p}_t + \eta_t (\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t)$ .
9:   Sample a constraint according to distribution  $\mathbf{p}$ .
10:  Calculate the stochastic gradient  $\hat{\nabla}_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t)$  and  $\hat{\nabla}_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t)$ .
11:  Randomly sample a constraint.
12:  Calculate the stochastic gradient  $\hat{\nabla}_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t)$ .
13:  Update  $m_{t,1}, m_{t,2}, m_{t,3}$ .
14:  Update  $\hat{m}_{t,1}, \hat{m}_{t,2}, \hat{m}_{t,3}$  and clip to  $[\rho, b]$ .
15:  Calculate  $A_{t,1}, A_{t,2}, A_{t,3}$ .
16: end while

```

randomly sample another element  $w_j$ , we can calculate the value of  $c_j$  and obtain the stochastic gradient of  $\mathcal{L}$  w.r.t  $\mathbf{p}$  as follows

$$\hat{\nabla}_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t) = de_j(c_j^2(\mathbf{w}_t, \boldsymbol{\lambda}_t; \mu^k) - \lambda p_j), \quad (6)$$

where  $e_j$  denotes a vector where its  $j$ -th element is 1 and other elements are 0. Since  $\mathbf{p}$  is related to the value of constraints, sampling constraint according to  $\mathbf{p}$  helps us find the most violating conditions.

To achieve a better performance, the momentum-based variance reduction method and adaptive method are also used. Specifically, we calculate the momentum-based gradient estimation w.r.t  $\mathbf{w}, \boldsymbol{\lambda}$

and  $\mathbf{p}$  as follows,

$$\begin{aligned} m_{t,1} &= \hat{\nabla}_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t) + (1 - a_{t+1,1}) \left( m_{t-1,1} - \hat{\nabla}_{\mathbf{w}} \mathcal{L}(\mathbf{w}_{t-1}, \boldsymbol{\lambda}_{t-1}, \mathbf{p}_{t-1}, \mu^k; \xi_t) \right) \\ m_{t,2} &= \hat{\nabla}_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t) + (1 - a_{t+1,2}) \left( m_{t-1,2} - \hat{\nabla}_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}_{t-1}, \boldsymbol{\lambda}_{t-1}, \mathbf{p}_{t-1}, \mu^k; \xi_t) \right) \\ m_{t,3} &= \hat{\nabla}_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t) + (1 - a_{t+1,3}) \left( m_{t-1,3} - \hat{\nabla}_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t-1}, \boldsymbol{\lambda}_{t-1}, \mathbf{p}_{t-1}, \mu^k; \xi_t) \right) \end{aligned}$$

Then, we calculate the adaptive matrix matrices  $A_{t,1}$ ,  $A_{t,2}$  and  $A_{t,3}$  for updating  $\mathbf{w}$ ,  $\boldsymbol{\lambda}$  and  $\mathbf{p}$ , respectively. Here, we present the calculation of adaptive matrix  $A_{t,1}$  as an example. Specifically, we calculate a second momentum-based estimation  $\hat{m}_{t,1} = \hat{a} \hat{\nabla}_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \boldsymbol{\lambda}_t, \mathbf{p}_t, \mu^k; \xi_t)^2 + (1 - \hat{a}) m_{t-1,1}$ . Then, we clip each element of  $\hat{m}_{t,1}$  into the range of  $[\rho, b]$  and obtain the adaptive matrix  $A_{t,1} = \text{diag} \left( \sqrt{\text{clip}(\hat{m}_{t,1}, \rho, b)} \right)$ . Note that we can use other method to calculate the adaptive matrices, such as AdaGrad-Norm (Ward et al., 2020), AMSGrad (Reddi et al., 2019), Adam<sup>+</sup> (Liu et al., 2020). Then, we can obtain the update rules as follows,

$$\tilde{\mathbf{w}}_{t+1} = \mathbf{w}_t - \gamma A_{t,1}^{-1} m_{t,1}, \quad \mathbf{w}_{t+1} = \mathbf{w}_t + \eta_t (\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t), \quad (7)$$

$$\tilde{\boldsymbol{\lambda}}_{t+1} = \boldsymbol{\lambda}_t - \gamma A_{t,2}^{-1} m_{t,2}, \quad \boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t + \eta_t (\tilde{\boldsymbol{\lambda}}_{t+1} - \boldsymbol{\lambda}_t), \quad (8)$$

$$\tilde{\mathbf{p}}_{t+1} = \mathcal{P}_{\Delta} (\mathbf{p}_t + \sigma A_{t,3}^{-1} m_{t,3}), \quad \mathbf{p}_{t+1} = \mathbf{p}_t + \eta_t (\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t), \quad (9)$$

where  $\gamma > 0$ ,  $\sigma > 0$ ,  $\eta_t > 0$  and  $\mathcal{P}_{\Delta}(\cdot)$  denotes the projection onto  $\Delta^d$ .

Once the following conditions are satisfied,

$$\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k)\|_2^2 \leq \epsilon_k^2, \quad \|\nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k)\|_2^2 \leq \epsilon_k^2, \quad \|c(\mathbf{w}, \boldsymbol{\lambda}; \mu^k)\|_2^2 \leq \epsilon_k^2, \quad (10)$$

where  $\nabla_{\mathbf{w}} \mathcal{L}$  and  $\nabla_{\boldsymbol{\lambda}} \mathcal{L}$  denote the full gradients  $\mathcal{L}$  w.r.t  $\mathbf{w}$  and  $\boldsymbol{\lambda}$ , we enlarge the penalty parameter  $\beta$ , and decrease the smooth parameter  $\mu^k$ .

The whole algorithm is presented in Algorithm 1 and 2. Note that instead of checking the conditions in each iteration of SCD, we check the conditions after several iterations to save time.

## 4 THEORETICAL ANALYSIS

In this section, we discuss the convergence performance of our proposed method (Detailed proofs are given in our appendix). Here, we give several assumptions which are widely used in convergence analysis.

**Assumption 1.** We have  $\mathbb{E}[\hat{\nabla}_{\mathbf{z}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)] = \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)$ ,  $\mathbb{E}[\hat{\nabla}_{\mathbf{p}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)] = \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)$ ,  $\mathbb{E}[\|\hat{\nabla}_{\mathbf{z}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k) - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2] \leq \mathbb{V}_{\mathbf{z}}^2$  and  $\mathbb{E}[\|\hat{\nabla}_{\mathbf{p}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k) - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2] \leq \mathbb{V}_{\mathbf{p}}^2$ .

**Assumption 2.** The function  $\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k)$  is  $\tau$ -strongly concave on  $\mathbf{p}$  for any give  $\mu^k > 0$ .

**Assumption 3.** The function  $\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k)$  has a  $L_1$ -Lipschitz gradient on  $(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p})$  for any give  $\mu^k > 0$ .

**Assumption 4.** The smoothing function  $\tilde{h}(\mathbf{w}, \mu^k)$  is twice continuously differentiable on  $\mathbf{w}$  for any  $\mu^k > 0$ .

**Assumption 5.**  $f$  is Lipschitz continuous and  $g$  is twice Lipschitz continuous w.r.t.  $\mathbf{w}$  and  $\boldsymbol{\lambda}$ .

We prove our Algorithm 2 can converge to the points satisfying conditions 10. Here, we give the definitions of  $\epsilon$ -stationary of the constrained problem 3 and minimax problem 4 and then show the relations between these definitions as follows,

**Definition 2.** ( $\epsilon$ -stationary point of the constrained optimization problem.)  $(\mathbf{w}^*, \boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*)$  is said to be the  $\epsilon$ -stationary point of the sub-problem (3) if the following conditions hold,  $\|\nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*; \mu^k) + \sum_{i=1}^d \alpha_i^* \nabla_{\mathbf{w}} c_i(\mathbf{w}^*, \boldsymbol{\lambda}^*)\|_2^2 \leq \epsilon_1^2$ ,  $\|\nabla_{\boldsymbol{\lambda}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*) + \sum_{i=1}^d \alpha_i^* \nabla_{\boldsymbol{\lambda}} c_i(\mathbf{w}^*, \boldsymbol{\lambda}^*; \mu^k)\|_2^2 \leq \epsilon_2^2$  and  $\sum_{i=1}^d c_i^2(\mathbf{w}^*, \boldsymbol{\lambda}^*; \mu^k) \leq \epsilon_3^2$ , where  $\boldsymbol{\alpha}$  denotes the lagrangian multipliers.

**Remark 1.** Let  $\alpha_j^* = p_j^* 2c_j(\mathbf{w}^*, \boldsymbol{\lambda}^*; \mu^k)$ . The conditions in Defintion 2 is equivalent to the tolerance conditions 10.

**Definition 3.** ( $\epsilon$ -stationary point of the mini-max problem.)  $(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mathbf{p}^*)$  is said to be the  $\epsilon$ -stationary point of the mini-max problem if it satisfies the conditions  $\|\nabla_{\mathbf{w}} \mathcal{L}\|_2^2 \leq \epsilon^2$ ,  $\|\nabla_{\boldsymbol{\lambda}} \mathcal{L}\|_2^2 \leq \epsilon^2$  and  $\|\nabla_{\mathbf{p}} \mathcal{L}\|_2^2 \leq \epsilon^2$ .

**Proposition 1.** If Assumptions 2 and 3 hold,  $(\mathbf{w}^*, \boldsymbol{\lambda}, \mathbf{p}^*)$  is the  $\epsilon$ -stationary point of the problem (4), then  $(\mathbf{w}^*, \boldsymbol{\lambda}^*)$  is the  $\epsilon$ -stationary point of the constrained problem 3.

According to Shi et al. (2022), the minimax problem 4 is equivalent to the following minimization problem:

$$\min_{\mathbf{w}, \boldsymbol{\lambda}} \left\{ H(\mathbf{w}, \boldsymbol{\lambda}) := \max_{\mathbf{p} \in \Delta^m} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}) = \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}^*(\mathbf{w}, \boldsymbol{\lambda})) \right\}, \quad (11)$$

where  $\mathbf{p}^*(\mathbf{w}, \boldsymbol{\lambda}) = \arg \max_{\mathbf{p}} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p})$ . Here, we give stationary point the minimization problem 11 and its relationship with Definition 3 as follows,

**Definition 4.** We call  $\mathbf{w}^*$  an  $\epsilon$ -stationary point of a differentiable function  $H(\mathbf{w})$ , if  $\|\nabla_{\mathbf{w}} H(\mathbf{w}^*, \boldsymbol{\lambda}^*)\|_2 \leq \epsilon$  and  $\|\nabla_{\boldsymbol{\lambda}} H(\mathbf{w}^*, \boldsymbol{\lambda}^*)\|_2 \leq \epsilon$ .

**Proposition 2.** Under Assumptions 3 and 2, if  $(\mathbf{w}', \boldsymbol{\lambda}')$  is the  $\epsilon$ -stationary point of  $H(\mathbf{w}, \boldsymbol{\lambda})$ , then  $(\mathbf{w}', \boldsymbol{\lambda}', \mathbf{p}')$  is the  $\epsilon$ -stationary point of  $\min_{\mathbf{w}, \boldsymbol{\lambda}} \max_{\mathbf{p} \in \Delta^d} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k)$  can be obtained. Conversely, if  $(\mathbf{w}', \boldsymbol{\lambda}', \mathbf{p}')$  is the  $\epsilon$ -stationary point of  $\min_{\mathbf{w}, \boldsymbol{\lambda}} \max_{\mathbf{p} \in \Delta^d} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k)$ , then a point  $(\mathbf{w}', \boldsymbol{\lambda}')$  is stationary point of  $H(\mathbf{w}, \boldsymbol{\lambda})$ .

**Remark 2.** According to Proposition 1 and Proposition 2, we have that once we find the  $\epsilon$ -stationary point in terms of Definition 4, then we can get the  $\epsilon$ -stationary point in terms of Definition 2. Therefore, we can obtain the points satisfying the tolerance conditions (10).

Before, we give the convergence result of our method, we present the lemma useful in our analysis. We have

**Lemma 1.** Under assumptions, let  $\mathbf{z} = [\mathbf{w}; \boldsymbol{\lambda}]$ , we have

$$\|\nabla H(\mathbf{z}) - m_{t,z}\|_2^2 \leq 2L_1^2 \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + 2\|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 \quad (12)$$

Then, we can define the following metric

$$\mathcal{M}_k = \frac{b^2}{\gamma^2} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 + 2L_1^2 \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + 2\|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2. \quad (13)$$

We have

$$\mathcal{M}_k \geq \frac{b^2}{\gamma^2} \|\mathbf{z}_t - \gamma A_{t,z}^{-1} m_{t,z} - \mathbf{z}_t\|_2^2 + \|\nabla H(\mathbf{z}) - m_{t,z}\|_2^2 \geq \|m_{t,z}\|_2^2 + \|\nabla H(\mathbf{z}) - m_{t,z}\|_2^2 \geq \frac{1}{2} \nabla H(\mathbf{z}_t)$$

If  $\mathcal{M}_k \rightarrow 0$ , we have  $\|\nabla H(\mathbf{z}_t)\|_2^2 \rightarrow 0$ . Thus, we can bound  $\mathcal{M}_k$  to find the stationary point of problem (11). Then, we give the convergence theorem in the following theorem.

**Theorem 1.** Assume Assumptions hold, if  $a_{t+1,1} = c_1 \eta_t^2, a_{t+1,2} = c_2 \eta_t^2, c_1 > \frac{5}{2} + \frac{2}{3e^3} \eta_b, 0 < \gamma \leq$

$$\frac{\sqrt{3}\tau\sigma\rho}{2\sqrt{12L_1^2\sigma^2\kappa^2 + 125L_1^2\kappa^2b^2}} \text{ and } 0 < \sigma \leq \min\left\{\frac{15}{12\mu}, \frac{1}{6L_1}\right\}, \text{ we have}$$

$$\frac{\gamma}{4\rho} \frac{1}{T} \sum_{t=1}^T \mathcal{M}_t \leq \frac{(\Theta_1 - \Theta^*)(m+T)^{1/3}}{Te} + \frac{\gamma(c_1^2 \mathbb{V}_z^2 + c_2^2 \mathbb{V}_p^2)(m+T)^{1/3}}{\rho T} \ln(m+T) \quad (14)$$

**Remark 3.** Theorem 1 demonstrate that with suitable setting, our method can converge to the points satisfying the conditions (10) at the rate of  $\tilde{\mathcal{O}}(T^{-2/3})$  if omitting log.

Then, we discuss the convergence performance of our whole algorithm. Define a new function

$$h_i^{\bar{\mathbf{w}}}(\mathbf{D}_i^T \mathbf{w}) := \begin{cases} h_i(\mathbf{D}_i^T \mathbf{w}) & i \notin \mathcal{I}_{\bar{\mathbf{w}}} \\ h_i(\mathbf{D}_i^T \bar{\mathbf{w}}) & i \in \mathcal{I}_{\bar{\mathbf{w}}} \end{cases}, \quad (15)$$

which is Lipschitz continuous at  $\mathbf{D}_i^T \bar{\mathbf{w}}, i = 1, 2, \dots, n$ . Then, we have  $h_{\bar{\mathbf{w}}}(\mathbf{w}) := (h_1^{\bar{\mathbf{w}}}(\mathbf{D}_1^T \mathbf{w}), h_2^{\bar{\mathbf{w}}}(\mathbf{D}_2^T \mathbf{w}), \dots, h_n^{\bar{\mathbf{w}}}(\mathbf{D}_n^T \mathbf{w}))$ , which has the same value as  $h(\mathbf{w})$  but opposite property.

For convenience, we define  $\phi_{\bar{w}}(\mathbf{w}) = \varphi(h_{\bar{w}}(\mathbf{w}))$  and  $\phi(\mathbf{w}) = \varphi(h(\mathbf{w}))$ . Besides, we define a vector set as follows,

$$\mathcal{V}_{\bar{w}} = \{\mathbf{v} : \mathbf{D}_i^T \mathbf{v} = \mathbf{0}, i \in \mathcal{I}_{\bar{w}}\}, \quad (16)$$

which means that  $\mathbf{v}$  is perpendicular to all column vectors in  $\mathbf{D}_i, i \in \mathcal{I}_{\bar{w}}$ . According to Bian & Chen (2017), the necessary condition of the non-Lipschitz lower level problem is

$$\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}})^T \mathbf{v} + \exp(\lambda_1) \phi^\circ(\mathbf{w}^*; \mathbf{v}) \geq 0, \quad (17)$$

for all  $\mathbf{v} \in \mathcal{V}_{\mathbf{w}^*}$ , where  $\phi^\circ(\mathbf{w}^*; \mathbf{v}) = \limsup_{\substack{\mathbf{w} \mapsto \mathbf{w}^* \\ t \downarrow 0}} \frac{\phi(\mathbf{w} + t\mathbf{v}) - \phi(\mathbf{w})}{t}$  denotes the Clarke generalized

directional derivative of  $\varphi(h(\mathbf{w}))$  at  $\mathbf{w}^*$ . Replacing the lower-level problem with above condition, we can obtain the following single-level problem,

$$\min_{\mathbf{w}, \boldsymbol{\lambda}} f(\mathbf{w}, \boldsymbol{\lambda}) \quad s.t. \quad c(\mathbf{w}, \boldsymbol{\lambda}) = \nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}})^T \mathbf{v} + \exp(\lambda_1) \phi^\circ(\mathbf{w}; \mathbf{v}) \geq 0 \quad (18)$$

for all  $\mathbf{v} \in \mathcal{V}_{\mathbf{w}^*}$ . For this new problem, we have the following theorem.

**Theorem 2.** *If  $(\mathbf{w}^*, \boldsymbol{\lambda}^*)$  satisfy the following conditions, then they are the stationary points of the problem (18).*

$$\nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*)^T \mathbf{v}_2 - (\mathbf{v}_2^T \nabla_{\mathbf{w}\mathbf{w}}^2 g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*) \mathbf{v}_1 + \exp(\lambda_1^*) \phi^{\circ\circ}(\mathbf{w}^*; \mathbf{v}_1, \mathbf{v}_2)) \xi^* \geq 0, \quad (19)$$

$$\nabla_{\boldsymbol{\lambda}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*)^T \mathbf{v}_3 - (\bar{\mathbf{v}}_3^T \nabla_{\mathbf{w}\bar{\boldsymbol{\lambda}}}^2 g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*) \mathbf{v}_1 + v_3^1 \exp(\lambda_1^*) \phi^\circ(\mathbf{w}^*; \mathbf{v}_1)) \xi^* \geq 0, \quad (20)$$

$$\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \exp(\lambda_1^*) \phi^\circ(\mathbf{w}^*; \mathbf{v}_1) \geq 0, \quad (21)$$

$$\xi^* (\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \lambda_1^* \phi^\circ(\mathbf{w}^*; \mathbf{v}_1)) = 0, \quad (22)$$

$$\xi^* \geq 0, \quad (23)$$

for all  $\mathbf{v}_1 \in \mathcal{V}_{\mathbf{w}^*}$ ,  $\mathbf{v}_2 \in \mathbb{R}^d$ ,  $\mathbf{v}_3 \in \mathbb{R}^m$ , where  $\mathbf{v}_3 = [v_3^1, \bar{\mathbf{v}}_3^T]^T$ ,  $\bar{\mathbf{v}}_3^T = [v_3^2, \dots, v_3^m]^T$  and  $\phi^{\circ\circ}(\mathbf{w}^*; \mathbf{v}_1, \mathbf{v}_2) = \limsup_{\substack{\mathbf{w} \mapsto \mathbf{w}^* \\ s \downarrow 0}} \frac{\phi^\circ(\mathbf{w} + \mathbf{v}_2 s; \mathbf{v}_1) - \phi^\circ(\mathbf{w}; \mathbf{v}_1)}{s}$ . In addition,  $(\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3)$  is direction vector used in calculating the Clarke directional derivative.

Then, we show with decreasing the smoothing parameter and tolerance parameters, our method can converge to stationary point defined in Theorem 2 in the following theorem.

**Theorem 3.** *Suppose  $\{\epsilon_k\}_{k=1}^\infty$  are positive and convergent ( $\lim_{k \rightarrow \infty} \epsilon_k = 0$ ) sequences,  $\{\mu^k\}_{k=1}^\infty$  is a positive and convergent ( $\lim_{k \rightarrow \infty} \mu^k = 0$ ) sequence. Then any limit point of the sequence points generated by SPNBO satisfies the conditions (19)-(23).*

Then, we show the relations between the conditions 19-23 and the original nonsmooth bilevel problem (1).

**Theorem 4.** *Assume the lower level problem in problem (1) is strongly convex. If we have  $(\mathbf{w}^*, \boldsymbol{\lambda}^*)$  and  $\xi^* \geq 0$  satisfying the conditions (19)-(23), then  $(\mathbf{w}^*, \boldsymbol{\lambda}^*)$  is the stationary point of the original nonsmooth bilevel problem.*

## 5 EXPERIMENTS

In this section, we conduct experiments to demonstrate the superiority of our method in terms of accuracy and efficiency.

### 5.1 EXPERIMENTAL SETUP

We summarize the baseline methods used in our experiments as follows.

1. **Penalty.** The method proposed in Mehra & Hamm (2019). It formulates the bi-level optimization problem as a one-level optimization problem, and then use penalty method to solve the new problem.

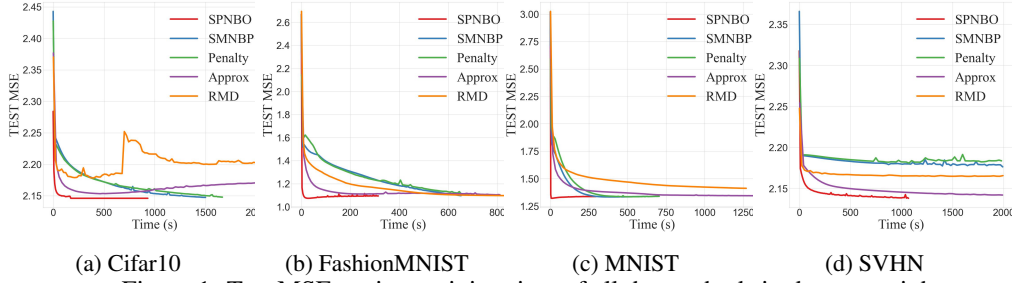


Figure 1: Test MSE against training time of all the methods in data re-weight.

2. **Approx.** The method proposed in Pedregosa (2016). It solves an additional linear problem to find the hypergradient to update the hyper-parameters.
3. **RMD.** The reverse method proposed in Franceschi et al. (2017). An additional loop is used to approximate the hypergradient.
4. **SMNBP.** The method proposed in Okuno et al. (2021). It uses the smoothing method to produce a sequence of smoothing lower-level functions and replaces them with the necessary condition. Then the penalty method is used to solve each single level problem.

We implement SMNBP, Penalty, Approx, RMD, and our method in Python. Since original Penalty, Approx and RMD are used for the smoothing problems, we use the smoothing function to approximate the lower-level problem. We fix the smoothing parameter at  $\mu = 0.0001$  in these methods. In SMNBP, for each given smoothing parameter, we solve the constrained problem using the Penalty method. For all these method, we use ADAM to update  $w$  and  $\lambda$  and choose the initial step size from  $\{0.1, 0.01, 0.001\}$ . For our method, we set  $\hat{a} = 0.9$  and other parameters are set according to Theorem 1. We choose the penalty parameter from  $\{0.1, 1, 10, 100\}$  for our method, SMNBP, and Penalty. We fix the inner iteration number  $T$  in Penalty, Approx, and RMD at 10 according to Mehra & Hamm (2019). We summarize the datasets used in our experiments in Table 2 and we divide all the datasets into three parts, *i.e.*, 40% for training, 40% for validation and 20% for testing. All the experiments are carried out 10 times on a PC with four 1080 Ti GPUs.

## 5.2 APPLICATIONS

**Data re-weight:** In this experiment, we evaluate the performance of all the methods in the application named data re-weight. In many real-world applications, the training set and testing set may have different distributions. To reduce the discrepancy between the two distributions, each data point will be given an additional importance weight, which is called data re-weight. In this application, we search the weight  $\lambda_i$  of each training data and the OSCAR regularization parameters  $\hat{\lambda}$  and  $\check{\lambda}$ . This problem can be formulated as

$$\begin{aligned}
 & \min_{\lambda} \frac{1}{N_{val}} \sum_{i=1}^{N_{val}} \ell(x_i^T w^*, y_i) \\
 & s.t. \ w^* \in \arg \min_w \sum_{i=1}^{N_{tr}} \frac{\exp(\lambda_i)}{\sum_j \exp(\lambda_j)} \ell(x_i^T w, y_i) + \exp(\hat{\lambda}) \|w\|_1 + \exp(\check{\lambda}) \sum_{i < j} \max\{w_{\mathcal{G}_i}, w_{\mathcal{G}_j}\},
 \end{aligned} \tag{24}$$

where  $N_{tr}$  and  $N_{val}$  denote the sizes of training set and validation set respectively.  $\{x_i, y_i\}$  denotes the data instance,  $\mathcal{G}_i$  denotes the group index. In this experiments, we set the group number equal to 10, and we use the mean squared loss  $\ell = (x_i^T w - y_i)^2$ .

**Training data poisoning:** In this experiment, we evaluate the performance of all the methods in training data poisoning. Assume we have pure training data  $\{x_i\}_{i=1}^{N_{tr}}$  with several poisoned points

Table 2: Datasets used in the experiments.

Datasets	Features	Samples	Classes
SVHN	$32 \times 32 \times 3$	73257	10
Cifar10	$84 \times 84 \times 3$	50000	10
MNIST	$28 \times 28 \times 1$	60000	10
FashionMNIST	$28 \times 28 \times 1$	60000	10



Table 3: The test mse of all the methods in data reweight. (Lower is better.)

Datasets	SPNBO	SMNBP	Penalty	Approx	RMD
Cifar10	<b>2.146</b> $\pm$ 0.006	2.147 $\pm$ 0.012	2.147 $\pm$ 0.011	2.171 $\pm$ 0.004	2.203 $\pm$ 0.012
MNIST	<b>1.338</b> $\pm$ 0.004	1.339 $\pm$ 0.006	1.340 $\pm$ 0.007	1.345 $\pm$ 0.010	1.412 $\pm$ 0.076
FashionMNIST	<b>1.091</b> $\pm$ 0.011	1.096 $\pm$ 0.020	1.100 $\pm$ 0.001	1.104 $\pm$ 0.009	1.097 $\pm$ 0.013
SVHN	<b>2.138</b> $\pm$ 0.004	2.176 $\pm$ 0.002	2.184 $\pm$ 0.006	2.142 $\pm$ 0.004	2.165 $\pm$ 0.002

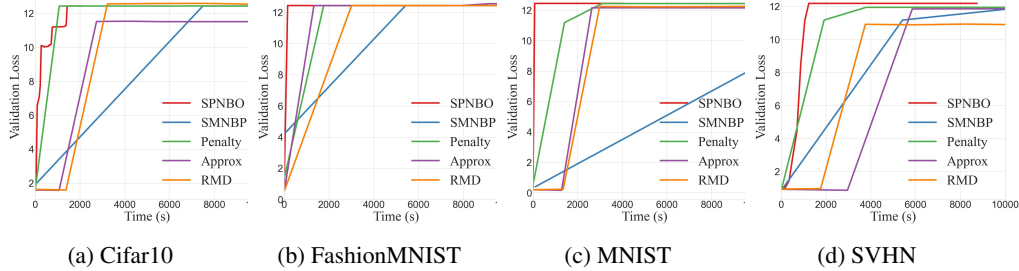


Figure 2: Validation loss versus training time of all the methods in training data poisoning.

$\{\lambda_j\}_{j=1}^{N_{poi}}$  assigned arbitrary labels. In this task, we search the poisoned data which can hurt the performance of the model trained from the clean data. This problem can be formulated as

$$\max_{\lambda} \frac{1}{N_{val}} \sum_{i=1}^{N_{val}} \ell(\theta(\mathbf{x}_i; \mathbf{w}^*), \mathbf{y}_i) \quad s.t. \quad \mathbf{w}^* \in \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{\mathbf{x}_i \in \mathcal{D}} \ell(\theta(\mathbf{x}_i; \mathbf{w}), \mathbf{y}_i) + \|\mathbf{w}\|_p^p,$$

where  $N = N_{tr} + N_{poi}$  and  $\mathcal{D}$  denote the dataset containing all the clean training data and poisoned data. In this experiment, we use Resnet18 as model. Besides, we add a  $p$ -norm ( $0 < p < 1$ ) regularization term in the lower-level problem to ensure that we can get a sparse model. In this experiment, we set  $p = 0.5$ . After solving the bilevel problem, we retrain the model on the clean data and poisoned data and then test the model.

Table 4: Test accuracy (%) of all the methods in training data poisoning (lower is better).

Data	Approx	RMD	Penalty	SMNBP	SPNBO
SVHN	50.79 $\pm$ 0.39	50.67 $\pm$ 0.27	50.67 $\pm$ 0.27	50.62 $\pm$ 0.29	<b>48.85</b> $\pm$ 0.57
Cifar10	82.91 $\pm$ 0.18	83.25 $\pm$ 0.11	82.29 $\pm$ 0.11	82.57 $\pm$ 0.11	<b>82.22</b> $\pm$ 0.28
FashionMNIST	96.09 $\pm$ 0.07	95.89 $\pm$ 0.31	95.87 $\pm$ 0.19	96.01 $\pm$ 0.22	<b>95.80</b> $\pm$ 0.20
MNIST	80.27 $\pm$ 0.25	77.63 $\pm$ 0.08	77.43 $\pm$ 0.54	77.50 $\pm$ 0.30	<b>77.22</b> $\pm$ 0.08

### 5.3 RESULTS AND DISCUSSION

All the results are presented in Tables 3, 4 and Figure 1, 2. From Table 3 and Table 4, we can find that our method has the similar results to other methods. From Figure 1 and Figure 2, we can find that our method is faster than other methods in most cases. This is because Approx and RMD need to solve the lower-level objective first and then need an additional loop to approximate the hypergradient which makes these methods have higher time complexity. Penalty and SMNBP need to use all the constraints in each updating step which is also time-consuming, when we use complex models (e.g., DNNs), Penalty and SMNBP suffer from high time complexity. However, our method uses the stochastic gradient method which makes it scalable to complicated models and does not need any intermediate steps to approximate the hypergradient. From all these results, we can conclude that our SPNBO is superior to other methods in terms of accuracy and efficiency.

## 6 CONCLUSION

In this paper, we proposed a new method, SPNBO, to solve the generalized non-smooth non-Lipschitz bi-level optimization problems by using the smoothing method and the penalty method. We also give the convergence analysis of our proposed method. The experimental results demonstrate the superiority of our method in terms of training time and accuracy.

## REFERENCES

- Jonathan F Bard. *Practical bilevel optimization: algorithms and applications*, volume 30. Springer Science & Business Media, 2013.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *The Journal of Machine Learning Research*, 13(1):281–305, 2012.
- James S Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pp. 2546–2554, 2011.
- Quentin Bertrand, Quentin Klopfenstein, Mathieu Blondel, Samuel Vaiter, Alexandre Gramfort, and Joseph Salmon. Implicit differentiation of lasso-type models for hyperparameter optimization. In *International Conference on Machine Learning*, pp. 810–821. PMLR, 2020.
- Dimitri P Bertsekas. On penalty and multiplier methods for constrained minimization. *SIAM Journal on Control and Optimization*, 14(2):216–235, 1976.
- Wei Bian and Xiaojun Chen. Optimality and complexity for constrained optimization problems with nonconvex regularization. *Mathematics of Operations Research*, 42(4):1063–1084, 2017.
- Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Nonconvex lagrangian-based optimization: monitoring schemes and global convergence. *Mathematics of Operations Research*, 43(4):1210–1232, 2018.
- Howard D Bondell and Brian J Reich. Simultaneous regression shrinkage, variable selection, and supervised clustering of predictors with oscar. *Biometrics*, 64(1):115–123, 2008.
- Kristian Bredies and Dirk A Lorenz. Linear convergence of iterative soft-thresholding. *Journal of Fourier Analysis and Applications*, 14(5-6):813–837, 2008.
- Alfred M Bruckstein, David L Donoho, and Michael Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. *SIAM review*, 51(1):34–81, 2009.
- X Chen, M Ng, and C Zhang. Nonconvex  $l_p$  regularization and box constrained model for image restoration. *IEEE Trans. Image Process*, 21.
- Xiaojun Chen, Lingfeng Niu, and Yaxiang Yuan. Optimality conditions and a smoothing trust region newton method for nonlipschitz optimization. *SIAM Journal on Optimization*, 23(3):1528–1552, 2013.
- Frank H Clarke. *Optimization and nonsmooth analysis*. SIAM, 1990.
- Benoît Colson, Patrice Marcotte, and Gilles Savard. An overview of bilevel optimization. *Annals of operations research*, 153(1):235–256, 2007.
- Andrew Cotter, Maya Gupta, and Jan Pfeifer. A light touch for heavily constrained sgd. In *Conference on Learning Theory*, pp. 729–771. PMLR, 2016.
- Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pp. 318–326, 2012.
- David L Donoho. De-noising by soft-thresholding. *IEEE transactions on information theory*, 41(3): 613–627, 1995.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.
- Matthias Feurer, Jost Tobias Springenberg, and Frank Hutter. Initializing bayesian hyperparameter optimization via meta-learning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and reverse gradient-based hyperparameter optimization. *arXiv preprint arXiv:1703.01785*, 2017.

- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning*, pp. 1568–1577. PMLR, 2018.
- Jordan Frecon, Saverio Salzo, and Massimiliano Pontil. Bilevel learning of the group lasso structure. *Advances in Neural Information Processing Systems*, 31:8301–8311, 2018.
- Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint arXiv:2007.05170*, 2020.
- Jian Huang, Joel L Horowitz, Shuangge Ma, et al. Asymptotic properties of bridge estimators in sparse high-dimensional regression models. *Annals of Statistics*, 36(2):587–613, 2008.
- A. Klein, S. Falkner, N. Mansur, and F. Hutter. Robo: A flexible and robust bayesian optimization framework in python. In *NIPS 2017 Bayesian Optimization Workshop*, December 2017.
- Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in neural information processing systems*, pp. 1008–1014. Citeseer, 2000.
- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020.
- Mingrui Liu, Wei Zhang, Francesco Orabona, and Tianbao Yang. Adam<sup>+</sup>: A stochastic method with adaptive variance reduction. *arXiv preprint arXiv:2011.11985*, 2020.
- Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *International Conference on Machine Learning*, pp. 2113–2122, 2015.
- Akshay Mehra and Jihun Hamm. Penalty method for inversion-free deep bilevel optimization. *arXiv preprint arXiv:1911.03432*, 2019.
- Harikrishna Narasimhan, Andrew Cotter, Yichen Zhou, Serena Wang, and Wenshuo Guo. Approximate heavily-constrained learning with lagrange multiplier models. *Advances in Neural Information Processing Systems*, 33:8693–8703, 2020.
- Yu Nesterov. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1): 127–152, 2005.
- Mila Nikolova, Michael K Ng, Shuqin Zhang, and Wai-Ki Ching. Efficient reconstruction of piecewise constant images using nonsmooth nonconvex minimization. *SIAM journal on Imaging Sciences*, 1(1):2–25, 2008.
- Takayuki Okuno, Akiko Takeda, Akihiro Kawana, and Motokazu Watanabe. On lp-hyperparameter learning via bilevel nonsmooth optimization. *J. Mach. Learn. Res.*, 22:245:1–245:47, 2021.
- Fabian Pedregosa. Hyperparameter optimization with approximate gradient. *arXiv preprint arXiv:1602.02355*, 2016.
- Aravind Rajeswaran, Chelsea Finn, Sham Kakade, and Sergey Levine. Meta-learning with implicit gradients. *arXiv preprint arXiv:1909.04630*, 2019.
- Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.
- Wanli Shi, Hongchang Gao, and Bin Gu. Gradient-free method for heavily constrained nonconvex optimization. In *International Conference on Machine Learning*, pp. 19935–19955. PMLR, 2022.
- Ankur Sinha, Pekka Malo, and Kalyanmoy Deb. Evolutionary algorithm for bilevel optimization using approximations of the lower level optimal solution mapping. *European Journal of Operational Research*, 257(2):395–411, 2017.
- Kevin Swersky, Jasper Snoek, and Ryan Prescott Adams. Freeze-thaw bayesian optimization. *arXiv preprint arXiv:1406.3896*, 2014.

Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes. *The Journal of Machine Learning Research*, 21(1):9047–9076, 2020.

Stephen Wright and Jorge Nocedal. Numerical optimization. *Springer Science*, 35(67-68):7, 1999.

Cun-Hui Zhang et al. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894–942, 2010.

## A CONVERGENCE ANALYSIS

### A.1 PROOF OF PROPOSITION 1

*Proof.* Since the  $(\mathbf{w}^*, \boldsymbol{\lambda}, \mathbf{p}^*)$  is the  $\epsilon$ -stationary point of  $\min_{\mathbf{w}} \max_{\mathbf{p} \in \Delta^d} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \mathbf{p}, \mu^k)$ , then we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*) + \beta \sum_{j=1}^d p_j^* 2c_j(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k) \nabla_{\mathbf{w}} c_j(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k)\|_2^2 \leq \epsilon^2. \quad (25)$$

Let  $\alpha_j^* = p_j^* 2c_j(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k)$  and  $\epsilon = \epsilon_1$ , we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*) + \sum_{j=1}^m \alpha_j^* \nabla_{\mathbf{w}} c_j(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k)\|_2^2 \leq \epsilon_1^2. \quad (26)$$

Then the first condition in Definition 2 is satisfied. Similarly, let  $\epsilon = \epsilon_2$ , we have

$$\|\nabla_{\boldsymbol{\lambda}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*) + \sum_{j=1}^m \alpha_j^* \nabla_{\boldsymbol{\lambda}} c_j(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k)\|_2^2 \leq \epsilon_2^2. \quad (27)$$

Since  $\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mathbf{p}^*)\|_2^2 \leq \epsilon^2$ , we have

$$\sum_{j=1}^d (\beta c_j^2(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k) - \tau p_j^*)^2 \leq \epsilon^2. \quad (28)$$

Using the inequality  $\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$  and  $0 \leq p_j \leq 1$ , we have

$$\begin{aligned} & \frac{1}{2} \beta^2 \sum_{j=1}^d c_j^4(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k) \\ & \leq \sum_{j=1}^d (\beta c_j^2(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k) - \tau p_j^*)^2 + \tau^2 \sum_{j=1}^d (p_j^*)^2 \\ & \leq \epsilon^2 + d\tau^2. \end{aligned} \quad (29)$$

Then, using  $(\frac{\sum_{i=0}^n a_i}{n})^2 \leq \frac{\sum_{i=0}^n a_i^2}{n}$ , we have

$$\sum_{j=1}^d c_j^2(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k) \leq \sqrt{d \sum_{j=1}^d c_j^4(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k)} \leq \sqrt{\frac{2d\epsilon^2 + 2d^2\tau^2}{\beta^2}}. \quad (30)$$

Let  $\sqrt{\frac{2d\epsilon^2 + 2d^2\tau^2}{\beta^2}} = \epsilon_3^2$ , we can obtain

$$\sum_{j=1}^m c_j^2(\mathbf{w}^*, \boldsymbol{\lambda}^*, \mu^k) \leq \epsilon_3^2. \quad (31)$$

Therefore, the second condition in Definition 2 is satisfied.

That completes the proof.  $\square$

## A.2 PROOF OF PROPOSITION 2

*Proof.* Let  $\mathbf{z} = [\mathbf{w}; \boldsymbol{\lambda}]$ . Assume that  $\mathbf{z}'$  satisfies the condition  $\|\nabla_{\mathbf{z}} H(\mathbf{z}')\|_2 \leq \epsilon$ . Then, solving the strongly concave problem  $\max_{\mathbf{p} \in \Delta^d} \mathcal{L}(\mathbf{z}', \mathbf{p}, \mu^k)$ , we can obtain a point  $\mathbf{p}'$  satisfying that

$$\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}', \mathbf{p}', \mu^k)\|_2 \leq \epsilon, \quad \|\mathbf{p}' - \mathbf{p}^*(\mathbf{z}')\|_2 \leq \epsilon. \quad (32)$$

In addition, we have

$$\begin{aligned} & \|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}', \mathbf{p}', \mu^k)\|_2 \\ & \leq \|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}', \mathbf{p}', \mu^k) - \nabla_{\mathbf{z}} H(\mathbf{z}')\|_2 + \|\nabla_{\mathbf{z}} H(\mathbf{z}')\|_2 \\ & = \|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}', \mathbf{p}', \mu^k) - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}', \mathbf{p}^*(\mathbf{z}'), \mu^k)\|_2 + \epsilon \\ & \leq L_1 \|\mathbf{p}' - \mathbf{p}^*(\mathbf{z}')\|_2 + \epsilon \\ & = \mathcal{O}(\epsilon). \end{aligned} \quad (33)$$

Thus,  $(\mathbf{z}', \boldsymbol{\lambda}')$  is the stationary point of  $\mathcal{L}$ .  $\square$

## A.3 PROOF OF LEMMA 1

*Proof.*

$$\begin{aligned} & \|\nabla H(\mathbf{z}_t) - m_{t,z}\|_2^2 \\ & = \|\nabla H(\mathbf{z}_t) - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) + \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 \\ & \leq 2\|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}^*(\mathbf{z}_t), \mu^k) - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2 + 2\|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 \\ & \leq 2L_1^2 \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + 2\|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 \end{aligned} \quad (34)$$

$\square$

## A.4 PROOF OF THEOREM 1

Before we discuss the convergence, we give several assumptions and definitions used in our analysis. All these assumptions are widely used in the analysis in Bolte et al. (2018); Mehra & Hamm (2019); Franceschi et al. (2017); Pedregosa (2016).

**Lemma 2.** (Lin et al. (2020)) Under Assumption 2 and Assumption 3, let  $\mathbf{z} = [\mathbf{w}^T, \boldsymbol{\lambda}^T]^T$ , the function  $H(\mathbf{z}) = \max_{\mathbf{p} \in \Delta^d} \mathcal{L}(\mathbf{z}, \mathbf{p}, \mu^k)$  has  $L_H = L_1 + \frac{L_1^2}{\tau}$ -Lipschitz continuous gradient.

**Lemma 3.** Under Assumption 2 and Assumption 3, let  $\mathbf{z} = [\mathbf{w}; \boldsymbol{\lambda}]$ , we have

$$H(\mathbf{z}_{t+1}) \leq H(\mathbf{z}_t) + \frac{\gamma \eta_t}{\rho} \|\nabla H(\mathbf{z}) - m_{t,z}\|_2^2 - \frac{\rho \eta_t}{2\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \quad (35)$$

*Proof.* According to Lemma 2, we have

$$\begin{aligned} & H(\mathbf{z}_{t+1}) \\ & \leq H(\mathbf{z}_t) + \langle \nabla H(\mathbf{z}), \mathbf{z}_{t+1} - \mathbf{z}_t \rangle + \frac{L_H}{2} \|\mathbf{z}_{t+1} - \mathbf{z}_t\|_2^2 \\ & = H(\mathbf{z}_t) + \langle \nabla H(\mathbf{z}), \eta_t (\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t) \rangle + \frac{L_H}{2} \|\eta_t (\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t)\|_2^2 \\ & = H(\mathbf{z}_t) + \langle \nabla H(\mathbf{z}), \eta_t (\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t) \rangle + \frac{L_H}{2} \|\eta_t (\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t)\|_2^2 \\ & = H(\mathbf{z}_t) + \eta_t \langle \nabla H(\mathbf{z}) - m_{t,z}, \tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t \rangle + \eta_t \langle m_{t,z}, \tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t \rangle + \frac{L_H \eta_t^2}{2} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \end{aligned} \quad (36)$$

According to the update rules of  $\mathbf{w}$  and  $\boldsymbol{\lambda}$ , we have  $\tilde{\mathbf{z}}_{t+1} = \mathbf{z}_t - \gamma A_{t,z}^{-1} m_{t,z} = \mathcal{P}(\mathbf{z}_t - \gamma A_{t,z}^{-1} m_{t,z})$ , where  $\tilde{\mathbf{z}}_{t+1} = [\tilde{\mathbf{w}}; \tilde{\boldsymbol{\lambda}}]$ ,  $m_{t,z} = [m_{t,1}, m_{t,2}]$  and  $A_{t,z} = \text{diag}([\hat{m}_{t,1}; \hat{m}_{t,2}] + \rho)$ . Then, we have

$$\langle \tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t + \gamma A_{t,z}^{-1} m_{t,z}, \mathbf{z} - \tilde{\mathbf{z}}_{t+1} \rangle \geq 0 \quad (37)$$

Set  $\mathbf{z} = \tilde{\mathbf{z}}$ , we have

$$\left\langle \frac{\gamma}{\rho} m_{t,z}, \mathbf{z}_t - \tilde{\mathbf{z}}_{t+1} \right\rangle \geq \langle \gamma A_{t,z}^{-1} m_{t,z}, \mathbf{z}_t - \tilde{\mathbf{z}}_{t+1} \rangle \geq \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \quad (38)$$

Thus, we have

$$\langle m_{t,z}, \tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t \rangle \leq -\frac{\rho}{\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \quad (39)$$

In addition, we have

$$\begin{aligned} & \langle \nabla H(\mathbf{z}) - m_{t,z}, \tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t \rangle \\ & \leq \|\nabla H(\mathbf{z}) - m_{t,z}\|_2 \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2 \\ & \leq \frac{\gamma}{\rho} \|\nabla H(\mathbf{z}) - m_{t,z}\|_2^2 + \frac{\rho}{4\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \end{aligned} \quad (40)$$

Thus, we have

$$\begin{aligned} & H(\mathbf{z}_{t+1}) \\ & \leq H(\mathbf{z}_t) + \eta_t \left( \frac{\gamma}{\rho} \|\nabla H(\mathbf{z}) - m_{t,z}\|_2^2 + \frac{\rho}{4\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \right) - \frac{\rho\eta_t}{\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 + \frac{L_H\eta_t^2}{2} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \\ & \leq H(\mathbf{z}_t) + \eta_t \frac{\gamma}{\rho} \|\nabla H(\mathbf{z}) - m_{t,z}\|_2^2 - \frac{\rho\eta_t}{2\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \end{aligned} \quad (41)$$

□

**Lemma 4.** Given  $0 < \eta_t \leq 1$ ,  $0 < \sigma \leq \frac{b}{6L_1}$ , for all  $t > 0$ , we have

$$\begin{aligned} \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 & \leq \left(1 - \frac{\eta_t\tau\sigma}{4b}\right) \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 - \frac{3\eta_t}{4} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 \\ & \quad + \frac{25\eta_t\sigma}{6\tau b} \|\nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + \frac{25\kappa^2\eta_t b}{6\tau\sigma} \|\mathbf{z}_t - \tilde{\mathbf{z}}_{t+1}\|_2^2 \end{aligned} \quad (42)$$

where  $\kappa = \frac{L_1}{\tau}$

*Proof.* Due to  $\mathcal{L}$  is stongly concave, we have

$$\begin{aligned} & \mathcal{L}(\mathbf{z}_t, \mathbf{p}, \mu^k) \\ & \leq \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) + \langle \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k), \mathbf{p} - \mathbf{p}_t \rangle - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 \\ & = \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) + \langle m_{t,p}, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle + \langle \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle \\ & \quad + \langle \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k), \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t \rangle - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 \end{aligned} \quad (43)$$

Since  $\mathcal{L}$  is  $L_1$  smooth, we have

$$-\frac{L_1}{2} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 \leq \mathcal{L}(\mathbf{z}_t, \tilde{\mathbf{p}}_{t+1}, \mu^k) - \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - \langle \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k), \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t \rangle \quad (44)$$

Thus, we have

$$\begin{aligned} & \mathcal{L}(\mathbf{z}_t, \mathbf{p}, \mu^k) \\ & \leq \mathcal{L}(\mathbf{z}_t, \tilde{\mathbf{p}}_{t+1}, \mu^k) + \frac{L_1}{2} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \langle m_{t,p}, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle + \langle \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle \\ & \quad - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 \end{aligned} \quad (45)$$

Since  $\tilde{\mathbf{p}}_{t+1} = \mathcal{P}_{\Delta}(\mathbf{p}_t + \sigma A_{t,3}^{-1} m_{t,p})$ , we have

$$\langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t - \sigma A_{t,3}^{-1} m_{t,p}, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle \geq 0 \quad (46)$$

Then, we obtain

$$\begin{aligned}
& \langle m_{t,p}, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle \\
& \leq \frac{b}{\sigma} \langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle \\
& = -\frac{b}{\sigma} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{b}{\sigma} \langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p} - \mathbf{p}_t \rangle
\end{aligned} \tag{47}$$

Thus, we have

$$\begin{aligned}
& \mathcal{L}(\mathbf{z}_t, \mathbf{p}, \mu^k) \\
& \leq \mathcal{L}(\mathbf{z}_t, \tilde{\mathbf{p}}_{t+1}, \mu^k) + \frac{L_1}{2} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 - \frac{b}{\sigma} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{b}{\sigma} \langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t \rangle \\
& \quad + \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p} - \tilde{\mathbf{p}}_{t+1} \rangle - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2
\end{aligned} \tag{48}$$

Set  $\mathbf{p} = \mathbf{p}^*(\mathbf{z}_t)$ , we have

$$\begin{aligned}
& \mathcal{L}(\mathbf{z}_t, \mathbf{p}^*(\mathbf{z}_t), \mu^k) \\
& \leq \mathcal{L}(\mathbf{z}_t, \tilde{\mathbf{p}}_{t+1}, \mu^k) + \frac{L_1}{2} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 - \frac{b}{\sigma} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{b}{\sigma} \langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t \rangle \\
& \quad + \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p}^*(\mathbf{z}_t) - \tilde{\mathbf{p}}_{t+1} \rangle - \frac{\tau}{2} \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2
\end{aligned} \tag{49}$$

Then, we have

$$\begin{aligned}
0 & \leq -\left(\frac{b}{\sigma} - \frac{L_1}{2}\right) \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{b}{\sigma} \langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t \rangle \\
& \quad + \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p}^*(\mathbf{z}_t) - \tilde{\mathbf{p}}_{t+1} \rangle - \frac{\tau}{2} \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2
\end{aligned} \tag{50}$$

By  $\mathbf{p}_{t+1} = \mathbf{p}_t + \eta_t(\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t)$ , we have

$$\begin{aligned}
& \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 \\
& = \|\mathbf{p}_t + \eta_t(\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t) - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 \\
& = \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 + 2\eta_t \langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t) \rangle + \eta_t^2 \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2
\end{aligned} \tag{51}$$

Then, we obtain

$$\langle \tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t \rangle = \frac{1}{2\eta_t} \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 + \frac{\eta_t}{2} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 - \frac{1}{2\eta_t} \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 \tag{52}$$

Also, we have

$$\begin{aligned}
& \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p}^*(\mathbf{z}_t) - \tilde{\mathbf{p}}_{t+1} \rangle \\
& = \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t \rangle + \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}, \mathbf{p}_t - \tilde{\mathbf{p}}_{t+1} \rangle \\
& \leq \frac{1}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + \frac{\tau}{4} \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + \frac{1}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + \frac{\tau}{4} \|\mathbf{p}_t - \tilde{\mathbf{p}}_{t+1}\|_2^2 \\
& = \frac{2}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + \frac{\tau}{4} \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}_t - \tilde{\mathbf{p}}_{t+1}\|_2^2
\end{aligned} \tag{53}$$

Combining above inequalities, we have

$$\begin{aligned}
& \frac{b}{2\eta_t\sigma} \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 \\
& \leq \left(\frac{b}{2\eta_t\sigma} - \frac{\tau}{4}\right) \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 + \left(\frac{b\eta_t}{2\sigma} + \frac{\tau}{4} + \frac{L_1}{2} - \frac{b}{\sigma}\right) \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{2}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 \\
& \leq \left(\frac{b}{2\eta_t\sigma} - \frac{\tau}{4}\right) \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 - \frac{3b}{8\sigma} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{2}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2.
\end{aligned} \tag{54}$$

where the last inequality is due to  $0 < \sigma \leq \frac{b}{6L_1}$ . Therefore, we have

$$\begin{aligned} & \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 \\ & \leq (1 - \frac{\eta_t \tau \sigma}{2b}) \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 - \frac{3\eta_t}{4} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{4\eta_t \sigma}{\tau b} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2. \end{aligned} \quad (55)$$

Next, we have

$$\begin{aligned} & \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 \\ & = \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t) + \mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 \\ & \leq \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 + \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 + 2\langle \mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t), \mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}^*(\mathbf{z}_{t+1}) \rangle \\ & \leq (1 + \frac{\eta_t \tau \sigma}{4b}) \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 + (1 + \frac{4b}{\eta_t \tau \sigma}) \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 \\ & \leq (1 + \frac{\eta_t \tau \sigma}{4b}) \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 + (1 + \frac{4b}{\eta_t \tau \sigma}) \kappa^2 \|\mathbf{z}_t - \mathbf{z}_{t+1}\|_2^2 \end{aligned} \quad (56)$$

Then, we have

$$\begin{aligned} & \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 \\ & \leq (1 + \frac{\eta_t \tau \sigma}{4b}) (1 - \frac{\eta_t \tau \sigma}{2b}) \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 - \frac{3\eta_t}{4} (1 + \frac{\eta_t \tau \sigma}{4b}) \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 \\ & \quad + \frac{4\eta_t \sigma}{\tau b} (1 + \frac{\eta_t \tau \sigma}{4b}) \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + (1 + \frac{4b}{\eta_t \tau \sigma}) \kappa^2 \|\mathbf{z}_t - \mathbf{z}_{t+1}\|_2^2 \\ & \leq (1 - \frac{\eta_t \tau \sigma}{4b}) \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 - \frac{3\eta_t}{4} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 \\ & \quad + \frac{25\eta_t \sigma}{6\tau b} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + \frac{25\kappa^2 \eta_t b}{6\tau \sigma} \|\mathbf{z}_t - \mathbf{z}_{t+1}\|_2^2 \end{aligned} \quad (57)$$

That complete the proof.  $\square$

**Lemma 5.** *Based on above assumptions, we have*

$$\begin{aligned} & \mathbb{E}[\|\mathbf{m}_{t+1,z} - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2^2] \\ & \leq (1 - a_{t+1,1}) \mathbb{E}[\|\mathbf{m}_{t,z} - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] + 2a_{t+1,1}^2 \mathbb{V}_z^2 + 4(1 - a_{t+1,1})^2 L_1^2 \eta_t^2 (\mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] + \mathbb{E}[\|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2]) \end{aligned} \quad (58)$$

$$\begin{aligned} & \mathbb{E}[\|\mathbf{m}_{t+1,p} - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2^2] \\ & \leq (1 - a_{t+1,2}) \mathbb{E}[\|\mathbf{m}_{t,p} - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] + 2a_{t+1,2}^2 \mathbb{V}_p^2 + 4(1 - a_{t+1,2})^2 L_1^2 \eta_t^2 (\mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] + \mathbb{E}[\|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2]) \end{aligned} \quad (59)$$



*Proof.* Since  $m_{t,z} = \hat{\nabla}_z \mathcal{L}(z_t, p_t, \mu^k; \xi_t) + (1 - a_t) (m_{t-1,z} - \hat{\nabla}_w \mathcal{L}(z_{t-1}, p_{t-1}, \mu^k; \xi_t))$ , we have

$$\begin{aligned}
& \mathbb{E}[\|m_{t+1,z} - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)\|_2^2] \\
&= \mathbb{E}[\|\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) + (1 - a_{t+1,1}) (m_{t,z} - \hat{\nabla}_w \mathcal{L}(z_t, p_t, \mu^k; \xi_{t+1})) - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)\|_2^2] \\
&= \mathbb{E}[\|(1 - a_{t+1,1}) (m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)) + a_{t+1,1} (\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)) \\
&\quad + (1 - a_{t+1,1}) (\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k) - \hat{\nabla}_z \mathcal{L}(z_t, p_t, \mu^k; \xi_{t+1}) + \nabla_z \mathcal{L}(z_t, p_t, \mu^k))\|_2^2] \\
&\leq (1 - a_{t+1,1})^2 \mathbb{E}[\|(m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k))\|_2^2] + \mathbb{E}[a_{t+1,1}^2 (\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)) \\
&\quad + (1 - a_{t+1,1}) (\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k) - \hat{\nabla}_z \mathcal{L}(z_t, p_t, \mu^k; \xi_{t+1}) + \nabla_z \mathcal{L}(z_t, p_t, \mu^k))\|_2^2] \\
&\leq (1 - a_{t+1,1})^2 \mathbb{E}[\|m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2a_{t+1,1}^2 \mathbb{E}[\|\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)\|_2^2] \\
&\quad + 2(1 - a_{t+1,1})^2 \mathbb{E}[\|\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k) - \hat{\nabla}_z \mathcal{L}(z_t, p_t, \mu^k; \xi_{t+1}) + \nabla_z \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] \\
&\leq (1 - a_{t+1,1})^2 \mathbb{E}[\|m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2a_{t+1,1}^2 \mathbb{V}_z^2 \\
&\quad + 2(1 - a_{t+1,1})^2 \mathbb{E}[\|\hat{\nabla}_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k; \xi_{t+1}) - \hat{\nabla}_z \mathcal{L}(z_t, p_t, \mu^k; \xi_{t+1})\|_2^2] \\
&\leq (1 - a_{t+1,1})^2 \mathbb{E}[\|m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2a_{t+1,1}^2 \mathbb{V}_z^2 + 4(1 - a_{t+1,1})^2 L_1^2 (\mathbb{E}[\|z_{t+1} - z_t\|_2^2] + \mathbb{E}[\|p_{t+1} - p_t\|_2^2]) \\
&\leq (1 - a_{t+1,1}) \mathbb{E}[\|m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2a_{t+1,1}^2 \mathbb{V}_z^2 + 4(1 - a_{t+1,1})^2 L_1^2 \eta_t^2 (\mathbb{E}[\|\tilde{z}_{t+1} - z_t\|_2^2] + \mathbb{E}[\|\tilde{p}_{t+1} - p_t\|_2^2]) \\
&\hspace{15em} (60)
\end{aligned}$$

Similarly, we have

$$\begin{aligned}
& \mathbb{E}[\|m_{t+1,p} - \nabla_p \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)\|_2^2] \\
&\leq (1 - a_{t+1,2}) \mathbb{E}[\|m_{t,p} - \nabla_p \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2a_{t+1,2}^2 \mathbb{V}_p^2 + 4(1 - a_{t+1,2})^2 L_1^2 \eta_t^2 (\mathbb{E}[\|\tilde{z}_{t+1} - z_t\|_2^2] + \mathbb{E}[\|\tilde{p}_{t+1} - p_t\|_2^2]) \\
&\hspace{15em} (61)
\end{aligned}$$

□

**Lemma 6.** Based on above results, we have  $a_{t+1,1} = c_1 \eta_t^2, a_{t+1,2} = c_2 \eta_t^2, c_1 > \frac{5}{2} + \frac{2}{3e^3} \eta_t$  and

$$c_2 > \frac{2}{3e^3} \eta_t + \frac{125L^2}{3\tau^2} \eta_t = \frac{e}{(m+t)^{1/3}},$$

$$\begin{aligned}
& \frac{1}{\eta_t} \mathbb{E}[\|m_{t+1,z} - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)\|_2^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|m_{t,z} - \nabla_z \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] \\
&\leq -\frac{5}{2} \eta_t \mathbb{E}[\|m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2c_1^2 \mathbb{V}_z^2 \eta_t^3 + 4L_1 \eta_t (\mathbb{E}[\|\tilde{z}_{t+1} - z_t\|_2^2] + \mathbb{E}[\|\tilde{p}_{t+1} - p_t\|_2^2]) \\
&\hspace{15em} (62)
\end{aligned}$$

$$\begin{aligned}
& \frac{1}{\eta_t} \mathbb{E}[\|m_{t+1,p} - \nabla_p \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)\|_2^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|m_{t,p} - \nabla_p \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] \\
&\leq -\frac{125L^2 \eta_t}{3\tau^2} \mathbb{E}[\|m_{t,p} - \nabla_p \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2c_2^2 \mathbb{V}_p^2 \eta_t^3 + 4L_1 \eta_t (\mathbb{E}[\|\tilde{z}_{t+1} - z_t\|_2^2] + \mathbb{E}[\|\tilde{p}_{t+1} - p_t\|_2^2]) \\
&\hspace{15em} (63)
\end{aligned}$$

*Proof.*

$$\begin{aligned}
& \frac{1}{\eta_t} \mathbb{E}[\|m_{t+1,z} - \nabla_z \mathcal{L}(z_{t+1}, p_{t+1}, \mu^k)\|_2^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|m_{t,z} - \nabla_z \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] \\
&\leq (\frac{1 - a_{t+1,1}}{\eta_t} - \frac{1}{\eta_{t-1}}) \mathbb{E}[\|m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + \frac{2a_{t+1,1}^2 \mathbb{V}_z^2}{\eta_t} + 4(1 - a_{t+1,1})^2 L_1 \eta_t (\mathbb{E}[\|\tilde{z}_{t+1} - z_t\|_2^2] + \mathbb{E}[\|\tilde{p}_{t+1} - p_t\|_2^2]) \\
&\leq (\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - c_1 \eta_t) \mathbb{E}[\|m_{t,z} - \nabla_w \mathcal{L}(z_t, p_t, \mu^k)\|_2^2] + 2c_1^2 \mathbb{V}_z^2 \eta_t^3 + 4L_1 \eta_t (\mathbb{E}[\|\tilde{z}_{t+1} - z_t\|_2^2] + \mathbb{E}[\|\tilde{p}_{t+1} - p_t\|_2^2]) \\
&\hspace{15em} (64)
\end{aligned}$$

Since  $\eta_t = \frac{e}{(m+t)^{1/3}}$  and  $(a+b)^{1/3} \leq a^{1/3} + ab^{-2/3}$ , we have

$$\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \leq \frac{2}{3e^3}\eta_t \quad (65)$$

Let  $c_1 > \frac{5}{2} + \frac{2}{3e^3}\eta_t$ , we have

$$\begin{aligned} & \frac{1}{\eta_t} \mathbb{E}[\|m_{t+1,z} - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|m_{t,z} - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] \\ & \leq -\frac{5}{2}\eta_t \mathbb{E}[\|m_{t,z} - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] + 2c_1^2 \mathbb{V}_z^2 \eta_t^3 + 4L_1^2 \eta_t (\mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] + \mathbb{E}[\|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2]) \end{aligned} \quad (66)$$

Similarly, Let  $c_2 > \frac{2}{3e^3}\eta_t + \frac{125L^2}{3\tau^2}$  we have

$$\begin{aligned} & \frac{1}{\eta_t} \mathbb{E}[\|m_{t+1,p} - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|m_{t,p} - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] \\ & \leq -\frac{125L^2\eta_t}{3\tau^2} \mathbb{E}[\|m_{t,p} - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] + 2c_2^2 \mathbb{V}_p^2 \eta_t^3 + 4L_1^2 \eta_t (\mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] + \mathbb{E}[\|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2]) \end{aligned} \quad (67)$$

□

Here, we give the proof of Theorem 1

*Proof.* According to lemma, we have

$$\begin{aligned} & \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 - \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 \\ & \leq -\frac{\eta_t \tau \sigma}{4b} \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 - \frac{3\eta_t}{4} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{25\eta_t \sigma}{6\tau b} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + \frac{25\kappa^2 \eta_t b}{6\tau \sigma} \|\mathbf{z}_t - \tilde{\mathbf{z}}_{t+1}\|_2^2 \end{aligned} \quad (68)$$

and

$$\begin{aligned} & H(\mathbf{z}_{t+1}) - H(\mathbf{z}_t) \\ & \leq \frac{\gamma \eta_t}{\rho} \|\nabla H(\mathbf{z}_t) - m_{t,z}\|_2^2 - \frac{\rho}{2\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \\ & \leq \frac{2L_1^2 \gamma \eta_t}{\rho} \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + \frac{2\gamma \eta_t}{\rho} \|\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 - \frac{\rho \eta_t}{2\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \end{aligned} \quad (69)$$

Define

$$\Theta_t = \mathbb{E}[H(\mathbf{z}_t) + \frac{10L_1^2 \gamma b}{\sigma \tau \rho} \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 + \frac{\gamma}{\rho \eta_{t-1}} (\|m_{t,z} - \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2 + \|m_{t,p} - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2)] \quad (70)$$

Then, we have

$$\begin{aligned}
& \Theta_{t+1} - \Theta_t \\
&= H(\mathbf{z}_{t+1}) - H(\mathbf{z}_t) + \frac{10L_1^2\gamma b}{\sigma\tau\rho} (\|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{z}_{t+1})\|_2^2 - \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2) \\
&\quad + \frac{\gamma}{\rho} \left( \frac{1}{\eta_t} \mathbb{E}[\|m_{t+1,z} - \nabla_{\mathbf{z}}\mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|m_{t,z} - \nabla_{\mathbf{z}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] \right. \\
&\quad \left. + \frac{1}{\eta_t} \mathbb{E}[\|m_{t+1,p} - \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_{t+1}, \mathbf{p}_{t+1}, \mu^k)\|_2^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|m_{t,p} - \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] \right) \\
&\leq \frac{2L_1^2\gamma\eta_t}{\rho} \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + \frac{2\gamma\eta_t}{\rho} \|\nabla_{\mathbf{z}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 - \frac{\rho\eta_t}{2\gamma} \|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2 \\
&\quad + \frac{10L_1^2\gamma b}{\sigma\tau\rho} \left( -\frac{\eta_t\tau\sigma}{4b} \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{z}_t)\|_2^2 - \frac{3\eta_t}{4} \|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{25\eta_t\sigma}{6\tau b} \|\nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,p}\|_2^2 + \frac{25\kappa^2\eta_t b}{6\tau\sigma} \|\mathbf{z}_t - \tilde{\mathbf{z}}_{t+1}\|_2^2 \right) \\
&\quad + \frac{\gamma}{\rho} \left( -\frac{5}{2} \eta_t \mathbb{E}[\|m_{t,z} - \nabla_{\mathbf{z}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] + 2c_1^2\mathbb{V}_z^2\eta_t^3 + 4L_1^2\eta_t (\mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] + \mathbb{E}[\|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2]) \right. \\
&\quad \left. - \frac{125L^2\eta_t}{3\tau^2} \mathbb{E}[\|m_{t,p} - \nabla_{\mathbf{p}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k)\|_2^2] + 2c_2^2\mathbb{V}_p^2\eta_t^3 + 4L_1^2\eta_t (\mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] + \mathbb{E}[\|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2]) \right) \\
&= \left( \frac{2L_1^2\gamma\eta_t}{\rho} - \frac{5L_1^2\gamma\eta_t}{2\rho} \right) \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 \\
&\quad + \left( \frac{2\gamma\eta_t}{\rho} - \frac{5\gamma\eta_t}{2\rho} \right) \|\nabla_{\mathbf{z}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 \\
&\quad + \left( -\frac{\rho\eta_t}{2\gamma} + \frac{10L_1^2\gamma b}{\sigma\tau\rho} \cdot \frac{25\kappa^2\eta_t b}{6\tau\sigma} + \frac{4L_1^2\eta_t\gamma}{\rho} \right) \mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] \\
&\quad + \left( -\frac{3\eta_t}{4} \cdot \frac{10L_1^2\gamma b}{\sigma\tau\rho} + \frac{4L_1^2\eta_t\gamma}{\rho} \right) \mathbb{E}[\|\tilde{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2] \\
&\quad + \frac{\gamma(c_1^2\mathbb{V}_z^2 + c_2^2\mathbb{V}_p^2)}{\rho} \eta_t^3 \\
&\leq -\frac{L_1^2\gamma\eta_t}{2\rho} \|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 - \frac{\gamma\eta_t}{2\rho} \|\nabla_{\mathbf{z}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 - \frac{\eta_t b^2}{4\rho\gamma} \mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2] + \frac{\gamma(c_1^2\mathbb{V}_z^2 + c_2^2\mathbb{V}_p^2)}{\rho} \eta_t^3
\end{aligned} \tag{71}$$

where the last inequality is due to  $0 < \gamma \leq \frac{\sqrt{3}\tau\sigma\rho}{2\sqrt{12L_1^2\sigma^2\kappa^2 + 125L_1^2\kappa^2b^2}}$  and  $0 < \sigma \leq \frac{15}{12\mu}$ .

Therefore, we have

$$\frac{\gamma\eta_t}{4\rho} (2L_1^2\|\mathbf{p}^*(\mathbf{z}_t) - \mathbf{p}_t\|_2^2 + 2\|\nabla_{\mathbf{z}}\mathcal{L}(\mathbf{z}_t, \mathbf{p}_t, \mu^k) - m_{t,z}\|_2^2 + \frac{b^2}{\gamma^2} \mathbb{E}[\|\tilde{\mathbf{z}}_{t+1} - \mathbf{z}_t\|_2^2]) \leq \Theta_t - \Theta_{t+1} + \frac{\gamma(c_1^2\mathbb{V}_z^2 + c_2^2\mathbb{V}_p^2)}{\rho} \eta_t^3 \tag{72}$$

Taking average over  $t = 1, \dots, T$ , we have  $\eta_t = \frac{e}{(m+t)^{1/3}}$

$$\begin{aligned}
& \frac{\gamma}{4\rho} \frac{1}{T} \sum_{t=1}^T \mathcal{M}_t \\
& \leq \frac{\Theta_1 - \Theta^*}{T\eta_K} + \frac{\gamma(c_1^2\mathbb{V}_z^2 + c_2^2\mathbb{V}_p^2)}{\rho T\eta_K} \sum_{t=1}^T \eta_t^3 \\
& \leq \frac{(\Theta_1 - \Theta^*)(m+T)^{1/3}}{Te} + \frac{\gamma(c_1^2\mathbb{V}_z^2 + c_2^2\mathbb{V}_p^2)(m+T)^{1/3}}{\rho Te} \int \frac{e}{m+t} dt \\
& = \frac{(\Theta_1 - \Theta^*)(m+T)^{1/3}}{Te} + \frac{\gamma(c_1^2\mathbb{V}_z^2 + c_2^2\mathbb{V}_p^2)(m+T)^{1/3}}{\rho T} \ln(m+T)
\end{aligned} \tag{73}$$

□

## B PROOF OF THEOREM 2

**Lemma 7.** *Under Assumption 5,  $c(\mathbf{w}, \boldsymbol{\lambda})$  is Lipschitz continuous on both  $\mathbf{w}$  and  $\boldsymbol{\lambda}$ .*

*Proof.* We have

$$\begin{aligned} & |c(\mathbf{w}, \boldsymbol{\lambda}) - c(\mathbf{w}', \boldsymbol{\lambda})| \\ &= |(\nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}})^T \mathbf{v} + \exp(\lambda_1) \phi^\circ(\mathbf{w}; \mathbf{v})) - (\nabla_{\mathbf{w}} g(\mathbf{w}', \bar{\boldsymbol{\lambda}})^T \mathbf{v} + \exp(\lambda_1) \phi^\circ(\mathbf{w}'; \mathbf{v}))| \\ &\leq |\nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}})^T \mathbf{v} - \nabla_{\mathbf{w}} g(\mathbf{w}', \bar{\boldsymbol{\lambda}})^T \mathbf{v}| + \exp(\lambda_1) |\phi^\circ(\mathbf{w}; \mathbf{v}) - \phi^\circ(\mathbf{w}'; \mathbf{v})| \end{aligned} \quad (74)$$

According to Proposition 2.1.1 in Clarke (1990) and  $g$  is  $L_g^w$ -Lipschitz continuous, where  $L_g^w$  is Lipschitz constant for  $\mathbf{w}$ , we have  $\nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}})^T \mathbf{v}$  is  $L_g^w$ -Lipschitz continuous. Besides, we have

$$\begin{aligned} & |\phi^\circ(\mathbf{w}; \mathbf{v}) - \phi^\circ(\mathbf{w}'; \mathbf{v})| \\ &= |\phi_{\bar{\mathbf{w}}}^\circ(\mathbf{w}; \mathbf{v}) - \phi_{\bar{\mathbf{w}}}^\circ(\mathbf{w}'; \mathbf{v})|. \end{aligned} \quad (75)$$

Since  $\phi_{\bar{\mathbf{w}}}^\circ$  is  $L_{\phi}^{\bar{\mathbf{w}}}$ -Lipschitz continuous according to Proposition 2.1.1 in Clarke (1990), where  $L_{\phi}^{\bar{\mathbf{w}}}$  is the Lipschitz constant for  $\mathbf{w}$ , we have  $\phi^\circ$  is  $L_{\phi}^{\bar{\mathbf{w}}}$ -Lipschitz continuous. Namely, we have

$$|\phi^\circ(\mathbf{w}; \mathbf{v}) - \phi^\circ(\mathbf{w}'; \mathbf{v})| \leq L_{\phi}^{\bar{\mathbf{w}}} \|\mathbf{w} - \mathbf{w}'\|_2 \quad (76)$$

Therefore, we have

$$|c(\mathbf{w}, \boldsymbol{\lambda}) - c(\mathbf{w}', \boldsymbol{\lambda})| \leq (L_g^w + L_{\phi}^{\bar{\mathbf{w}}}) \|\mathbf{w} - \mathbf{w}'\|_2, \quad (77)$$

which means that  $c(\mathbf{w}, \boldsymbol{\lambda})$  is Lipschitz continuous on  $\mathbf{w}$ .

Then we prove  $c$  is Lipschitz continuous on  $\boldsymbol{\lambda}$ . We have

$$\begin{aligned} & |c(\mathbf{w}, \boldsymbol{\lambda}) - c(\mathbf{w}, \boldsymbol{\lambda}')| \\ &= |(\nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}})^T \mathbf{v} + \exp(\lambda_1) \phi^\circ(\mathbf{w}; \mathbf{v})) - (\nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}}')^T \mathbf{v}_1 + \exp(\lambda'_1) \phi^\circ(\mathbf{w}; \mathbf{v}_1))| \\ &= |\nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}})^T \mathbf{v} - \nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}}')^T \mathbf{v}| + |\exp(\lambda_1) - \exp(\lambda'_1)| \cdot |\phi^\circ(\mathbf{w}; \mathbf{v})| \end{aligned} \quad (78)$$

According to the Proposition 2.1.1 in Clarke (1990), we have

$$|\phi^\circ(\mathbf{w}; \mathbf{v})| = |\phi_{\bar{\mathbf{w}}}^\circ(\mathbf{w}; \mathbf{v})| \leq L_{\phi}^{\bar{\mathbf{w}}} \|\mathbf{v}\|_2 \quad (79)$$

Thus, we obtain

$$|c(\mathbf{w}, \boldsymbol{\lambda}) - c(\mathbf{w}, \boldsymbol{\lambda}')| \leq C \|\boldsymbol{\lambda} - \boldsymbol{\lambda}'\|_2. \quad (80)$$

where  $C$  is a constant. That completes the proof.  $\square$

Then, we prove Theorem 2

*Proof.* Define

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \xi) = f(\mathbf{w}, \boldsymbol{\lambda}) - c(\mathbf{w}, \boldsymbol{\lambda})^T \xi \quad (81)$$

where  $\xi \geq 0$ . If  $(\mathbf{w}^*, \boldsymbol{\lambda})$  is the stationary point of problem 18, we have

$$\mathbf{0} \in \partial_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*, \boldsymbol{\lambda}^*, \xi^*) \quad (82)$$

$$\mathbf{0} \in \partial_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}^*, \boldsymbol{\lambda}^*, \xi^*) \quad (83)$$

$$\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \lambda_1^* \phi^\circ(\mathbf{w}^*; \mathbf{v}_1) \geq 0 \quad (84)$$

$$\xi^* (\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \lambda_1^* \phi^\circ(\mathbf{w}^*; \mathbf{v}_1)) = 0 \quad (85)$$

$$\xi^* \geq 0 \quad (86)$$

Since  $c(\mathbf{w}, \boldsymbol{\lambda})$  and  $f(\mathbf{w}, \boldsymbol{\lambda})$  are Lipschitz continuous, we have

$$\begin{aligned} & |\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \xi) - \mathcal{L}(\mathbf{w}', \boldsymbol{\lambda}, \xi)| \\ &\leq |f(\mathbf{w}, \boldsymbol{\lambda}) - f(\mathbf{w}', \boldsymbol{\lambda})| + \xi |c(\mathbf{w}, \boldsymbol{\lambda}) - c(\mathbf{w}', \boldsymbol{\lambda})| \\ &\leq (L_f^w + L_g^w + L_{\phi}^{\bar{\mathbf{w}}}) \|\mathbf{w} - \mathbf{w}'\|_2, \end{aligned} \quad (87)$$

where  $L_f^w$  is the Lipschitz constant of  $f$  on  $w$ . This means  $\mathcal{L}(w, \lambda, \xi)$  is Lipschitz continuous. Thus, condition 82 is equivalent to

$$\mathcal{L}_w^\circ(w^*, \lambda^*, \xi^*; v_2) \geq 0, \forall v_2 \in \mathbb{R}^d \quad (88)$$

where  $\mathcal{L}_w^\circ(w^*, \lambda^*, \xi^*; v_2)$  is the generalized clark directional derivative at  $w^*$  defined as follows

$$\mathcal{L}_w^\circ(w^*, \lambda^*, \xi^*; v_2) = \limsup_{\substack{w \mapsto w^*, \\ s \downarrow 0}} \frac{\mathcal{L}(w + sv_2, \lambda^*, \xi^*) - \mathcal{L}(w, \lambda^*, \xi^*)}{s}. \quad (89)$$

Then, for all  $v_2$ , we have,

$$\begin{aligned} & \mathcal{L}_w^\circ(w^*, \lambda^*, \xi^*; v_2) \\ &= \limsup_{\substack{w \mapsto w^*, \\ s \downarrow 0}} \frac{\mathcal{L}(w + sv_2, \lambda^*, \xi^*) - \mathcal{L}(w, \lambda^*, \xi^*)}{s} \\ &= \limsup_{\substack{w \mapsto w^*, \\ s \downarrow 0}} \frac{f(w + sv_2, \lambda^*) - f(w, \lambda^*, \xi^*)}{s} \\ &\quad - \xi^* \limsup_{\substack{w \mapsto w^*, \\ s \downarrow 0}} \frac{\nabla_w g(w + sv_2, \lambda^*) v_1 - \nabla_w g(w, \lambda^*) v_1}{s} \\ &\quad - \exp(\lambda_1^*) \xi^* \limsup_{\substack{w \mapsto w^*, \\ s \downarrow 0}} \frac{\phi^\circ(w + sv_2; v_1) - \phi^\circ(w^*; v_1)}{s} \\ &= \nabla_w f(w^*, \lambda^*)^T v_2 - (v_2^T \nabla_{ww}^2 g(w^*, \bar{\lambda}^*) v_1 + \exp(\lambda_1^*) \phi^{\circ\circ}(w^*; v_1, v_2)) \xi^* \\ &\geq 0 \end{aligned} \quad (90)$$

where

$$\phi^{\circ\circ}(w^*; v_1, v_2) = \limsup_{\substack{w \mapsto w^*, \\ s \downarrow 0}} \frac{\phi^\circ(w + sv_2; v_1) - \phi^\circ(w; v_1)}{s} \quad (91)$$

Using the same method, we can obtain condition (83) is equivalent to

$$\nabla_\lambda f(w^*, \lambda^*)^T v_3 - (\bar{v}_3 \nabla_{w\bar{\lambda}}^2 g(w^*, \bar{\lambda}^*) v_1 + v_3^1 \exp(\lambda_1^*) \phi^\circ(w^*; v_1)) \xi^* \geq 0 \quad (92)$$

for all  $v_3 \in \mathbb{R}^d$ .

That completes the proof.  $\square$

### B.1 PROOF OF THEOREM 3

*Proof.* In our method, we solve the following smoothed single level problem for each given  $\mu^k$

$$\begin{aligned} & \min_{w, \lambda} f(w, \lambda) \\ & \text{s.t. } c(w, \lambda; \mu^k) := \nabla_w g(w, \bar{\lambda}) + \lambda_1 \nabla_w \varphi(\tilde{h}(w, \mu^k)) = \mathbf{0}, \end{aligned} \quad (93)$$

where  $\nabla_w \varphi(\tilde{h}(w, \mu^k)) = \varphi'(z)_{h(w, \mu^k)} \nabla_w h(w, \mu^k)$  and  $\mu^k$  is the smoothing parameter.

Let  $w^k, \lambda^k, p^k$  be the point return by the random constraint method, according to the tolerance condition, we have

$$\|c(w^*, \lambda^*; \mu^k)\|_2^2 \leq \epsilon_k^2 \quad (94)$$

Multiplying  $\|v_1\|_2^2$ , where  $v_1 \in \mathcal{V}_{w^*}$ , on both sides of above equality, we have

$$\|c(w^k, \lambda^k; \mu^k)\|_2 \|v_1\|_2 \leq \epsilon_k \|v_1\|_2 \quad (95)$$

According to  $|ab| \leq \|a\|_2 \|b\|_2$ , we have

$$|c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k)^T \mathbf{v}_1| \leq \epsilon_k \|\mathbf{v}_1\|_2. \quad (96)$$

Thus, we have

$$c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k)^T \mathbf{v}_1 \geq -\epsilon_k \|\mathbf{v}_1\|_2. \quad (97)$$

Then, taking the limit on both sides (i.e.,  $k \rightarrow \infty$ ), we have

$$\lim_{k \rightarrow \infty} c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k)^T \mathbf{v}_1 \geq 0 \quad (98)$$

Then, according to the Theorem 2 in Bian & Chen (2017), we have

$$\lim_{k \rightarrow \infty} c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k)^T \mathbf{v}_1 = \nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \lambda_1^* \phi^\circ(\mathbf{w}^*; \mathbf{v}_1) \quad (99)$$

Thus, we have

$$\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \lambda_1^* \phi^\circ(\mathbf{w}^*; \mathbf{v}_1) \geq 0 \quad (100)$$

Since we have  $\alpha_j^k = 2p_j^k c_j(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k)$  and  $\lim_{k \rightarrow \infty} \|c_j(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k)\|_2^2 = 0$ , we have  $\boldsymbol{\alpha}^* = \lim_{k \rightarrow \infty} \boldsymbol{\alpha}^k = \mathbf{0}$ . Therefore, we have  $\boldsymbol{\alpha}^* \in \mathcal{V}_{\mathbf{w}^*}$  and

$$c(\mathbf{w}^*, \boldsymbol{\lambda}^*; \mu^k)^T \boldsymbol{\alpha}^* = 0. \quad (101)$$

Let  $\boldsymbol{\alpha}^* = \xi^* \mathbf{v}_1 \in \mathcal{V}_{\mathbf{w}^*}$  and  $\xi^* = 0$ , we can obtain

$$\xi^* (\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \lambda_1^* \phi^\circ(\mathbf{w}^*; \mathbf{v}_1)) = 0 \quad (102)$$

We also have

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}^k, \boldsymbol{\lambda}^k) - \nabla_{\mathbf{w}} c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k) \boldsymbol{\alpha}^k\|_2^2 \leq \epsilon_k^2 \quad (103)$$

Assume that we have a vector  $\mathbf{v}_2$ . Multiplying  $\|\mathbf{v}_2\|_2^2$  on the both side of the above inequality, we have

$$\|\nabla_{\mathbf{w}} f(\mathbf{w}^k, \boldsymbol{\lambda}^k) - \nabla_{\mathbf{w}} c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k) \boldsymbol{\alpha}^k\|_2 \|\mathbf{v}_2\|_2 \leq \epsilon_k \|\mathbf{v}_2\|_2. \quad (104)$$

According to  $|\langle \mathbf{a}, \mathbf{b} \rangle| \leq \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$ , we have

$$|\langle \nabla_{\mathbf{w}} f(\mathbf{w}^k, \boldsymbol{\lambda}^k) - \nabla_{\mathbf{w}} c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k) \boldsymbol{\alpha}^k, \mathbf{v}_2 \rangle| \leq \epsilon_k \|\mathbf{v}_2\|_2 \quad (105)$$

Obviously, we have

$$\langle \nabla_{\mathbf{w}} f(\mathbf{w}^k, \boldsymbol{\lambda}^k) - \nabla_{\mathbf{w}} c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k) \boldsymbol{\alpha}^k, \mathbf{v}_2 \rangle \geq -\epsilon_k \|\mathbf{v}_2\|_2 \quad (106)$$

Taking the limit on both sides, we can obtain

$$\lim_{k \rightarrow \infty} \langle \nabla_{\mathbf{w}} f(\mathbf{w}^k, \boldsymbol{\lambda}^k) - \nabla_{\mathbf{w}} c(\mathbf{w}^k, \boldsymbol{\lambda}^k; \mu^k) \boldsymbol{\alpha}^k, \mathbf{v}_2 \rangle \geq 0 \quad (107)$$

Namely, we have

$$\langle \nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*) - \nabla_{\mathbf{w}} c(\mathbf{w}^*, \boldsymbol{\lambda}^*; \mu^k) \boldsymbol{\alpha}^*, \mathbf{v}_2 \rangle \geq 0 \quad (108)$$

In addition, we have

$$\begin{aligned} & \langle \nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*) - \nabla_{\mathbf{w}} c(\mathbf{w}^*, \boldsymbol{\lambda}^*; \mu^k) \boldsymbol{\alpha}^*, \mathbf{v}_2 \rangle \\ &= \langle \nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*) - \xi^* \nabla_{\mathbf{w}} (\nabla_{\mathbf{w}} g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*)^T \mathbf{v}_1 + \lambda_1^* \phi^\circ(\mathbf{w}^*; \mathbf{v}_1)), \mathbf{v}_2 \rangle \\ &= \nabla_{\mathbf{w}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*)^T \mathbf{v}_2 - \xi^* \mathbf{v}_2^T \nabla_{\mathbf{w}}^2 g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*) \mathbf{v}_1 - \lambda_1^* \phi^{\circ\circ}(\mathbf{w}^*; \mathbf{v}_1, \mathbf{v}_2) \\ &\geq 0 \end{aligned} \quad (109)$$

where

$$\phi^{\circ\circ}(\mathbf{w}^*; \mathbf{v}_1, \mathbf{v}_2) = \limsup_{\substack{\mathbf{w} \mapsto \mathbf{w}^*, \\ s \downarrow 0,}} \frac{\phi^\circ(\mathbf{w} + \mathbf{v}_2 s; \mathbf{v}_1) - \phi^\circ(\mathbf{w}; \mathbf{v}_1)}{s} \quad (110)$$

By using the same method, we can obtain

$$\nabla_{\boldsymbol{\lambda}} f(\mathbf{w}^*, \boldsymbol{\lambda}^*)^T \mathbf{v}_3 - \left( \bar{v}_3 \nabla_{\mathbf{w} \bar{\boldsymbol{\lambda}}}^2 g(\mathbf{w}^*, \bar{\boldsymbol{\lambda}}^*) \mathbf{v}_1 + v_3^1 e^{\lambda_1^*} \phi^\circ(\mathbf{w}^*; \mathbf{v}_1) \right) \xi^* \geq 0 \quad (111)$$

for all  $\mathbf{v}_3$ .

That completes the proof.  $\square$

Table 5: The test mse of our method in data reweight with different penalty parameter  $\beta$ .

Datasets	0.1	1	10	100
Cifar10	$2.186 \pm 0.014$	$2.157 \pm 0.016$	$2.148 \pm 0.007$	$2.146 \pm 0.006$
MNIST	$1.339 \pm 0.009$	$1.338 \pm 0.007$	$1.341 \pm 0.009$	$1.338 \pm 0.004$
FashionMNIST	$1.091 \pm 0.011$	$1.105 \pm 0.011$	$1.090 \pm 0.012$	$1.094 \pm 0.008$
SVHN	$2.140 \pm 0.003$	$2.141 \pm 0.002$	$2.138 \pm 0.004$	$2.154 \pm 0.043$

## C PROOF OF THEOREM 4

Here we prove Theorem 4

*Proof.* Since the lower level problem in problem (1) is convex, finding the optimal solution of the lower level problem is equivalent to find the point which satisfies the condition

$$\mathbf{0} \in \nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}}) + \exp(\lambda_1) \partial \varphi(h(\mathbf{w})) \quad (112)$$

This is equivalent to find the point which satisfies the condition

$$\nabla_{\mathbf{w}} g(\mathbf{w}, \bar{\boldsymbol{\lambda}})^T \mathbf{v} + \exp(\lambda_1) \phi^\circ(\mathbf{w}; \mathbf{v}) \geq 0, \forall \mathbf{v} \in \mathcal{V}_{\mathbf{w}} \quad (113)$$

Therefore,  $(\mathbf{w}^*, \boldsymbol{\lambda}^*)$  is the optimal solution of the single level problem (18). Thus,  $(\mathbf{w}^*, \boldsymbol{\lambda}^*)$  satisfies the conditions 19-23.  $\square$

## D IMPACT OF PENALTY PARAMETER

In this section, we evaluate the impact of different initial values for penalty parameter  $\beta$ . We fix other parameters and choose  $\beta$  from  $\{0.1, 1, 10, 100\}$ . We present the results in Table 5 and Table 6. From all these results, we can find that using a large penalty parameter usually leads to better results. This is because using a large penalty parameter means that we can impose greater penalties onto the conditions that violate the constraints, so as to ensure better convergence performance.

Table 6: Test accuracy (%) of our method in training data poisoning with different penalty parameter.

Data	0.1	1	10	100
SVHN	$48.90 \pm 0.37$	$49.05 \pm 0.23$	$48.95 \pm 0.24$	$48.85 \pm 0.57$
Cifar10	$83.20 \pm 0.22$	$82.76 \pm 0.16$	$82.32 \pm 0.09$	$82.22 \pm 0.28$
FashionMNIST	$96.02 \pm 0.12$	$95.99 \pm 0.18$	$96.06 \pm 0.16$	$95.80 \pm 0.20$
MNIST	$79.31 \pm 0.13$	$77.32 \pm 0.04$	$78.08 \pm 0.21$	$77.22 \pm 0.08$